

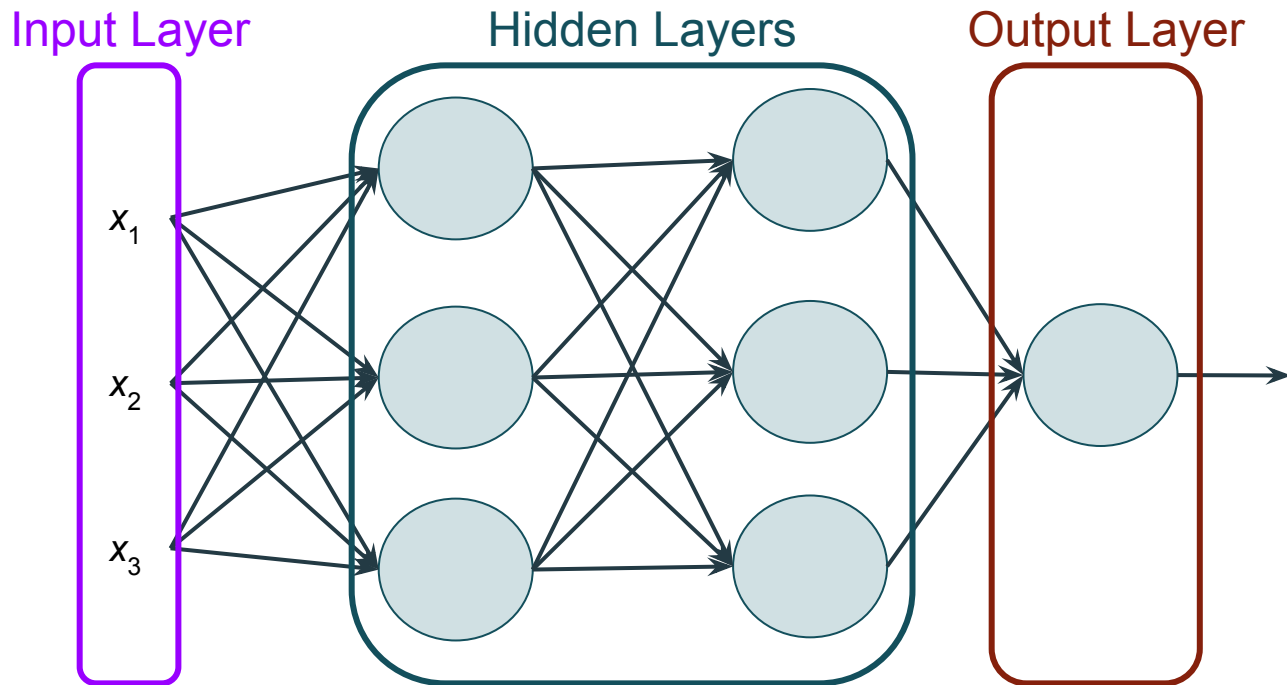
Advanced Neural Networks



CS344
Deep Learning



Deep Learning

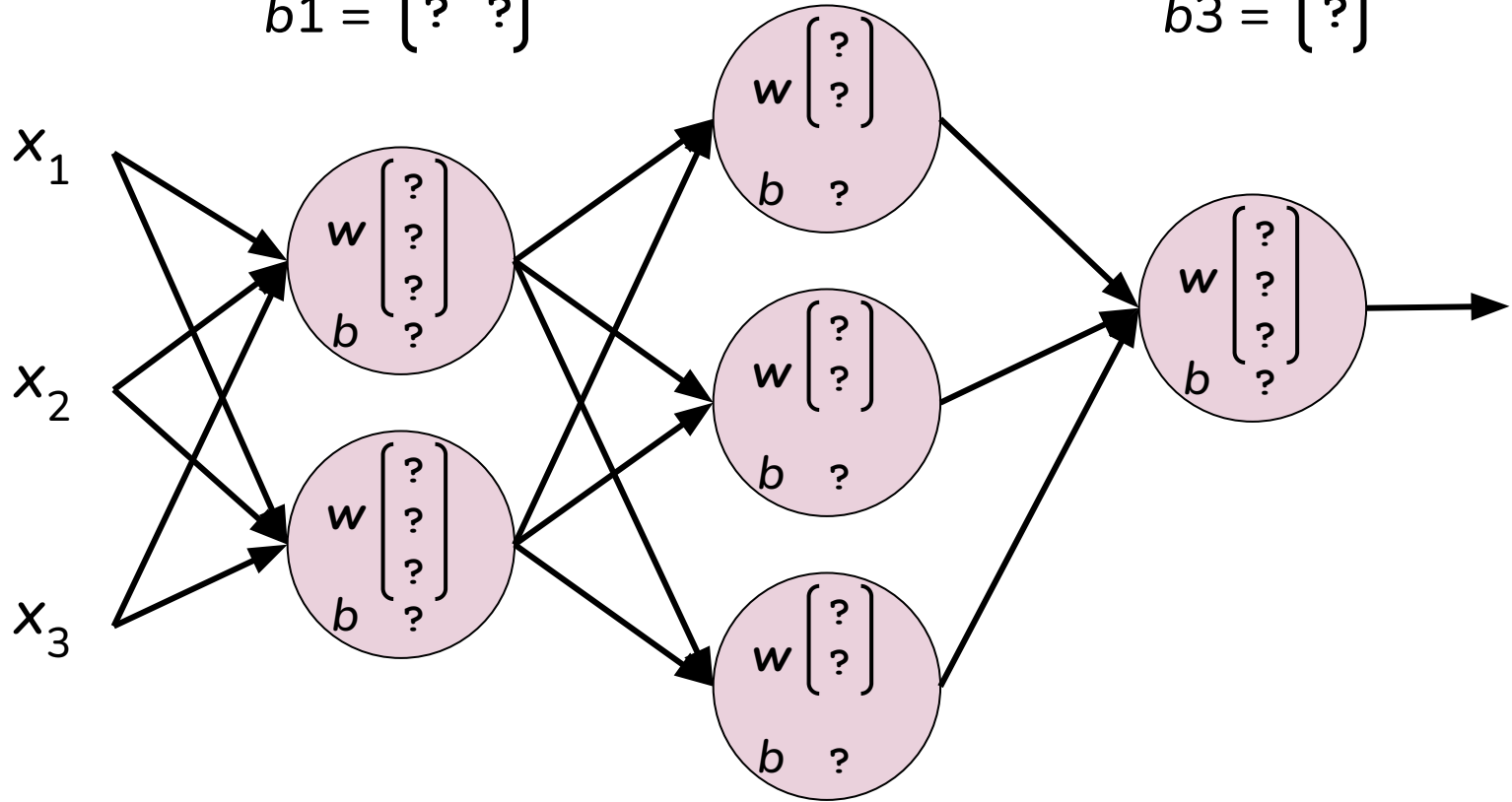


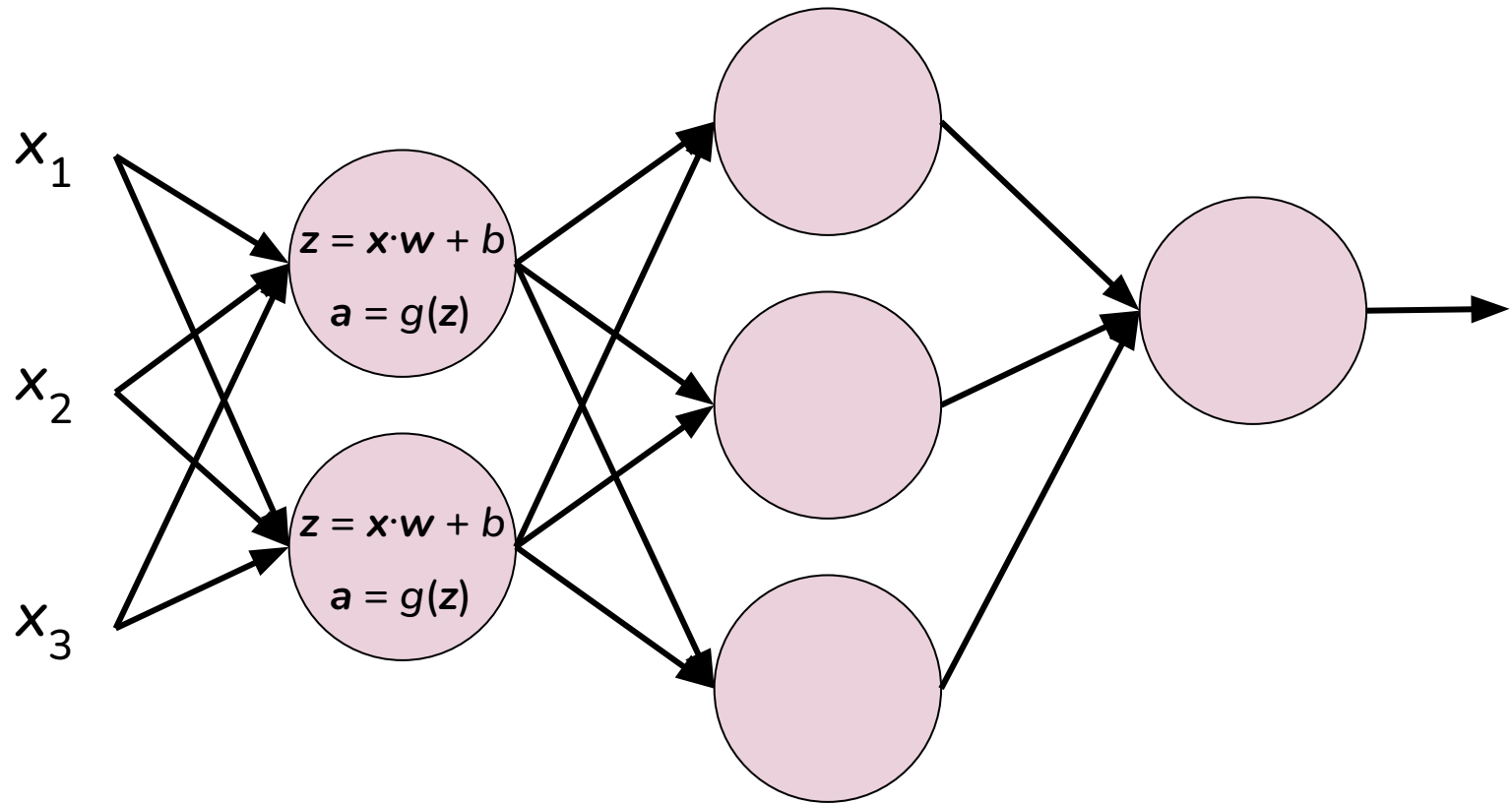
Hidden layers enable NNs to learn their own non-linear features!

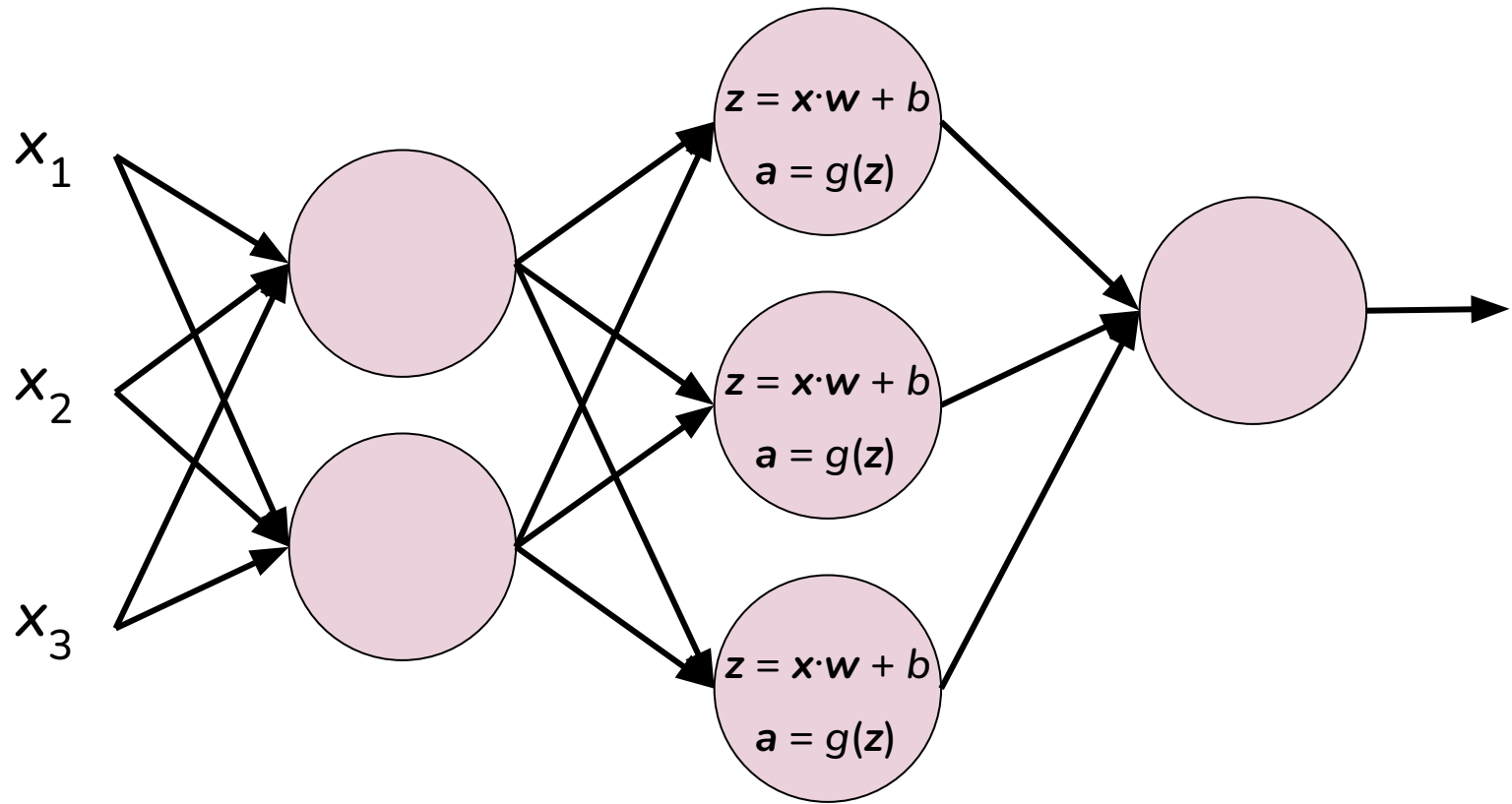
$$W1 = \begin{pmatrix} ? & ? \\ ? & ? \\ ? & ? \end{pmatrix}$$
$$b1 = \begin{pmatrix} ? & ? \end{pmatrix}$$

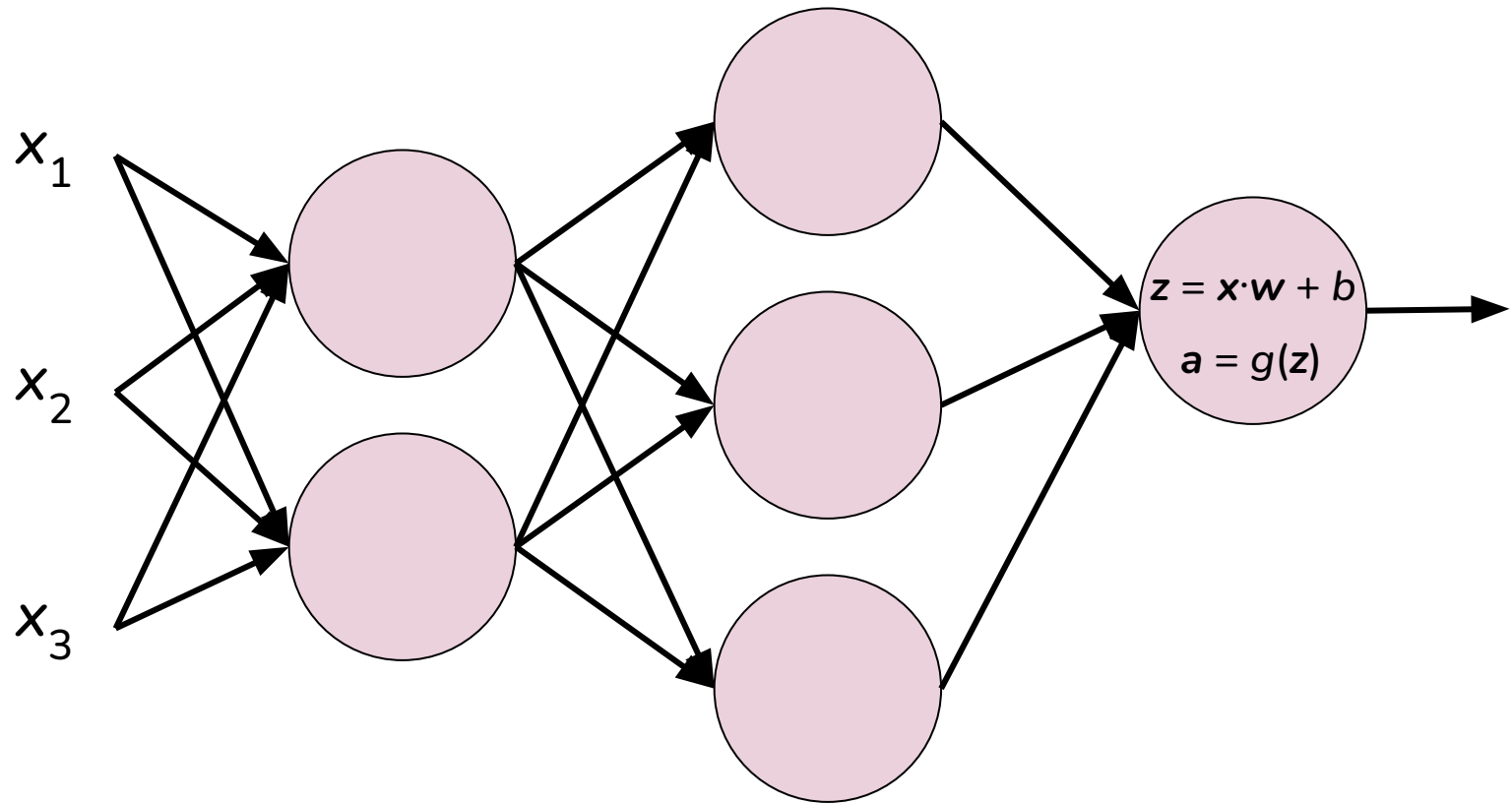
$$W2 = \begin{pmatrix} ? & ? & ? \\ ? & ? & ? \end{pmatrix}$$
$$b2 = \begin{pmatrix} ? & ? & ? \end{pmatrix}$$

$$W3 = \begin{pmatrix} ? \\ ? \\ ? \end{pmatrix}$$
$$b3 = \begin{pmatrix} ? \end{pmatrix}$$



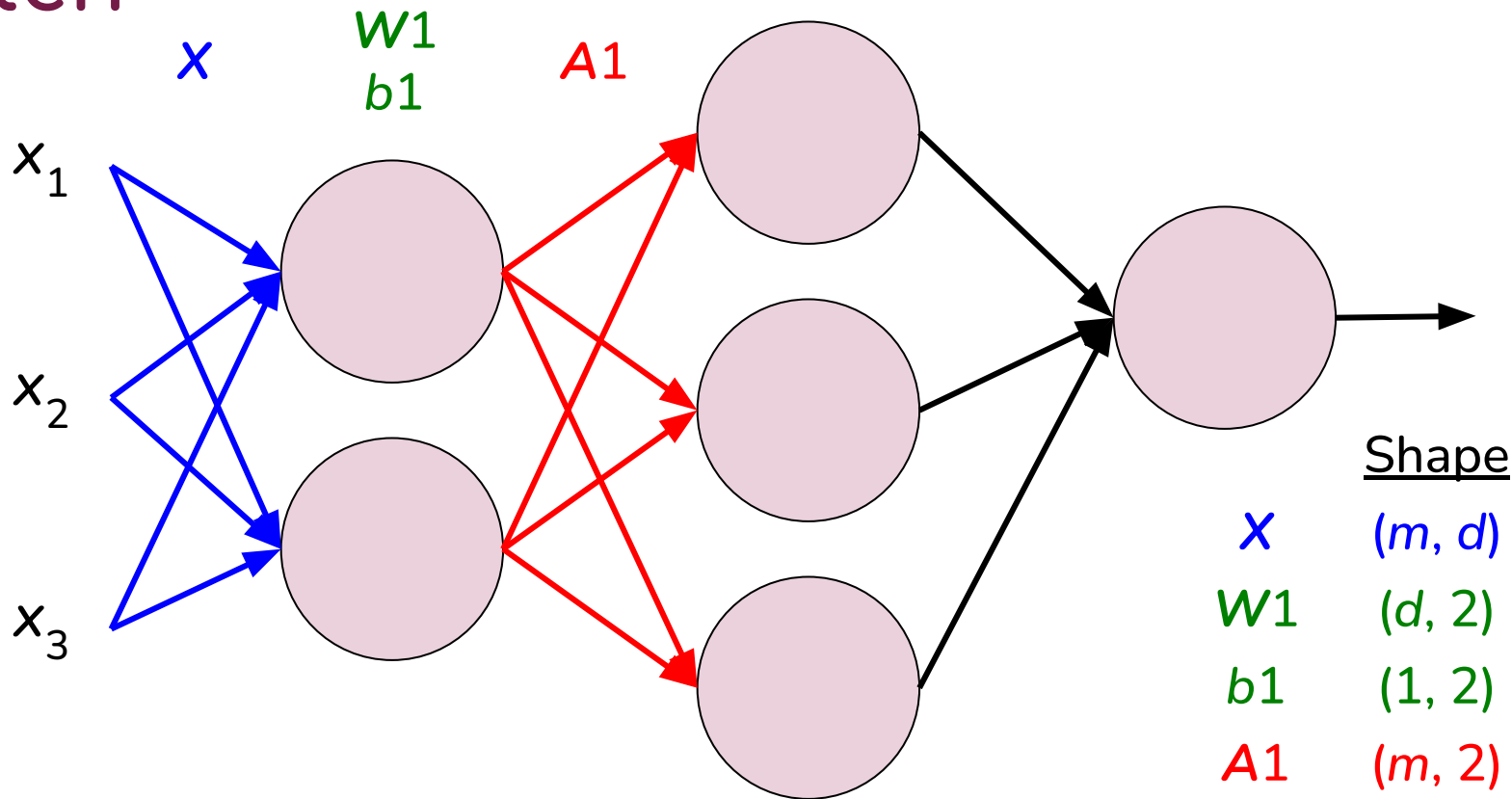






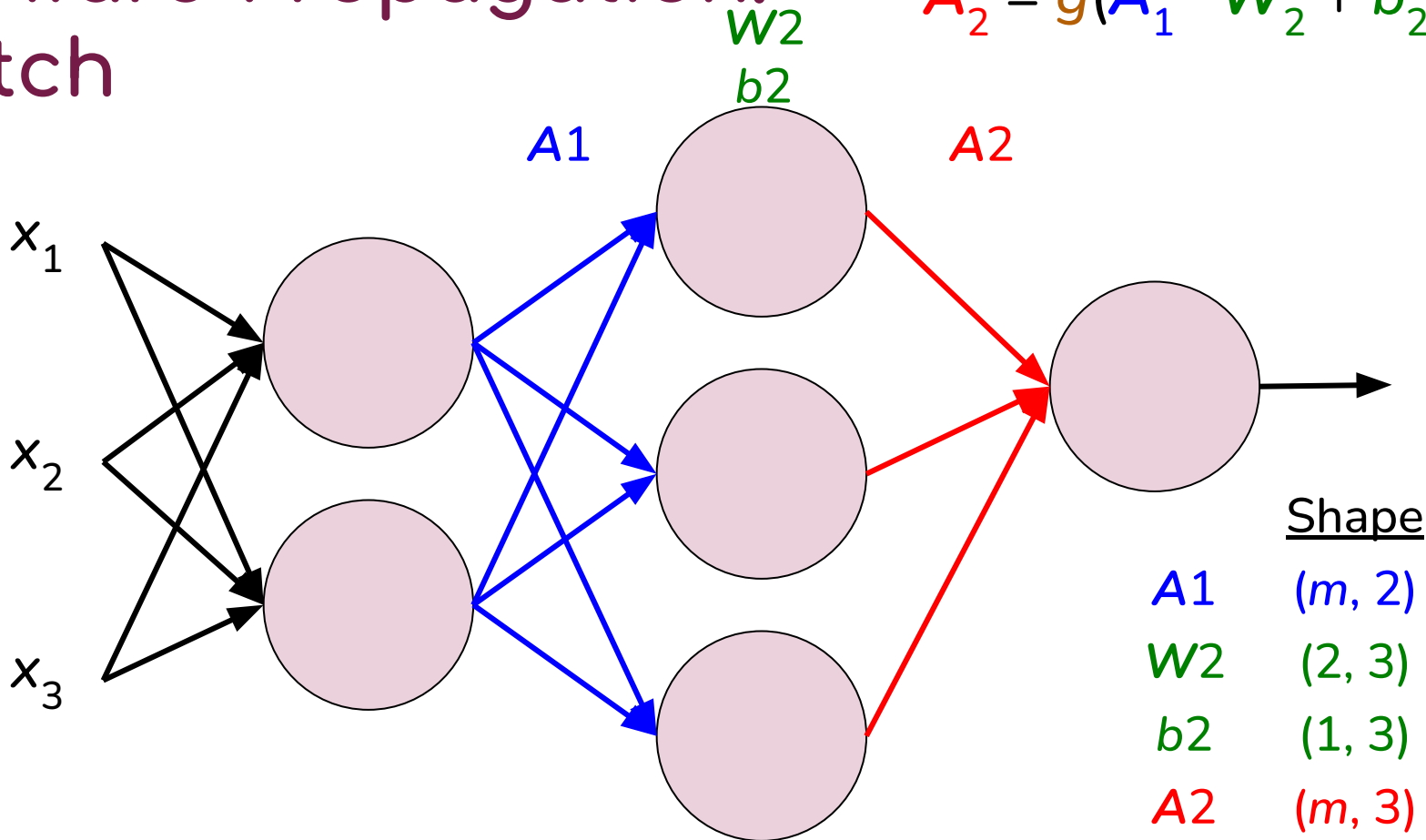
Forward Propagation: Batch

$$A_1 = g(X \cdot W_1 + b_1)$$



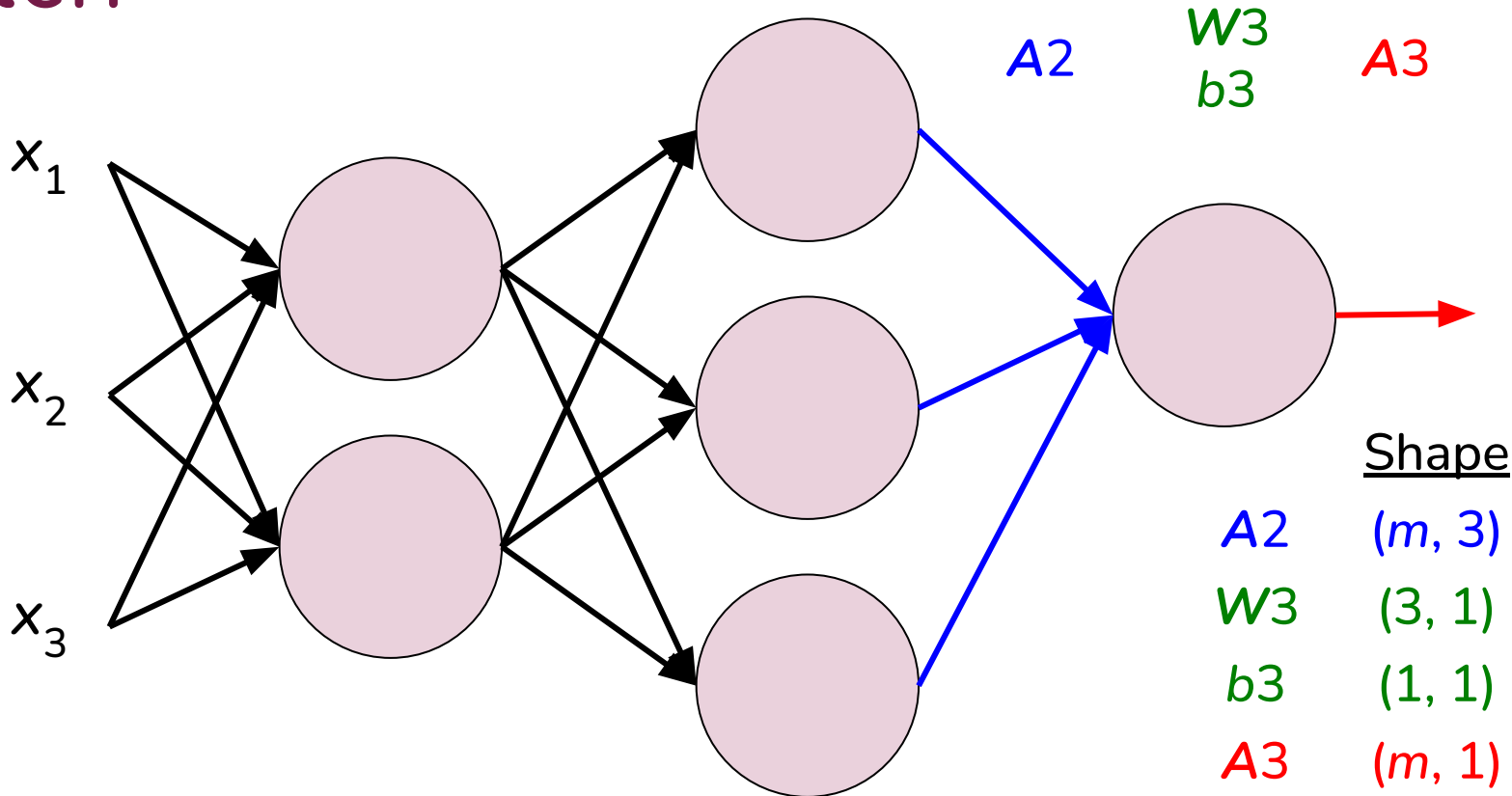
Forward Propagation: Batch

$$A_2 = g(A_1 \cdot W_2 + b_2)$$



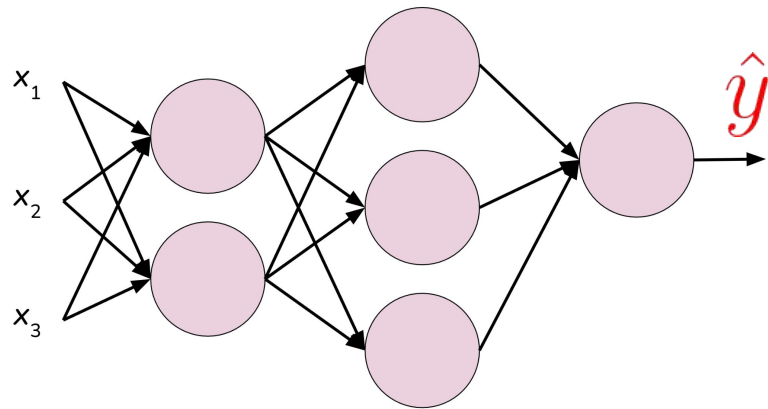
Forward Propagation: Batch

$$A_3 = g(A_2 \cdot W_3 + b_3)$$



Loss and Cost

The *loss function*, L , quantifies the error, i.e., how far our prediction \hat{y} is from the true label y



$$L = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

The *cost function*, J , is the average loss (error) over all data points

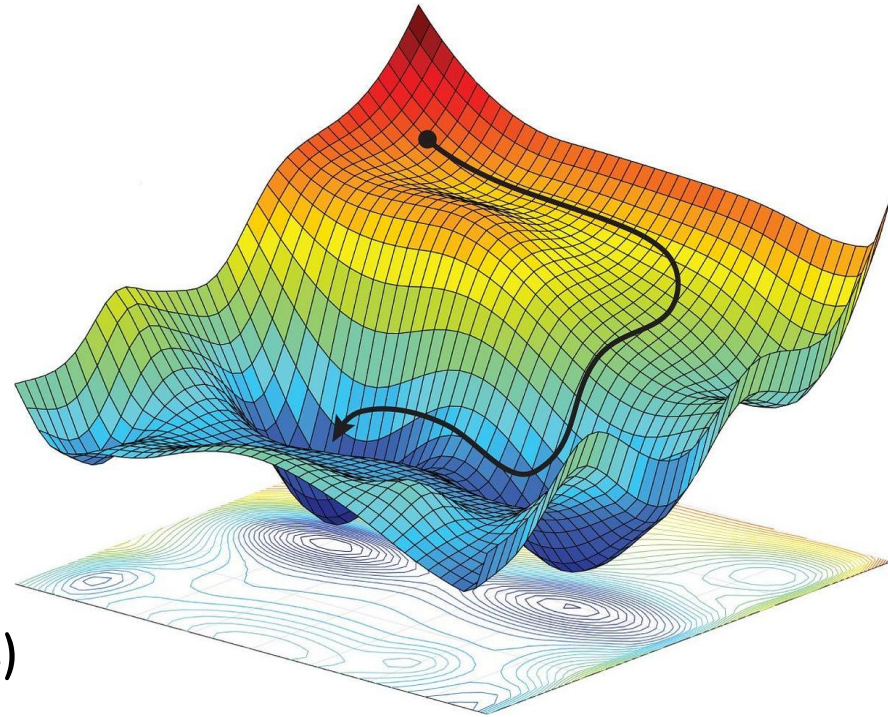
$$J = \frac{1}{m} \sum_{i=1}^m L = \frac{1}{m} \sum_{i=1}^m -y_i \log(\hat{y}_i) - (1 - y_i) \log(1 - \hat{y}_i)$$

Training

We want to find parameters (W 's and b 's) that minimize the cost, J

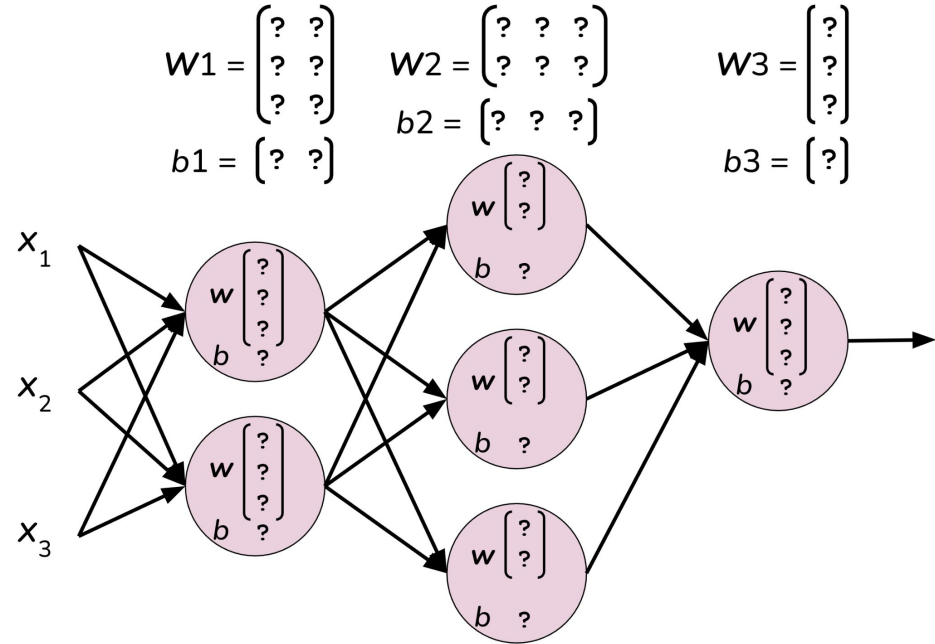
Gradient Descent Algorithm

- ❖ Initialize parameters (W 's and b 's)
- ❖ Repeat until converge:
 - Update parameters (W 's and b 's) to reduce the cost, J



Random Initialization

- ❖ If units in the same layer *start* with the same parameters then they will *end* with the same parameters
- ❖ There's no point in having repetitive units, i.e., multiple units in a layer with the same parameter values
- ❖ Thus, we initialize the parameters *randomly* so they start (and end) with different values



Gradient Descent

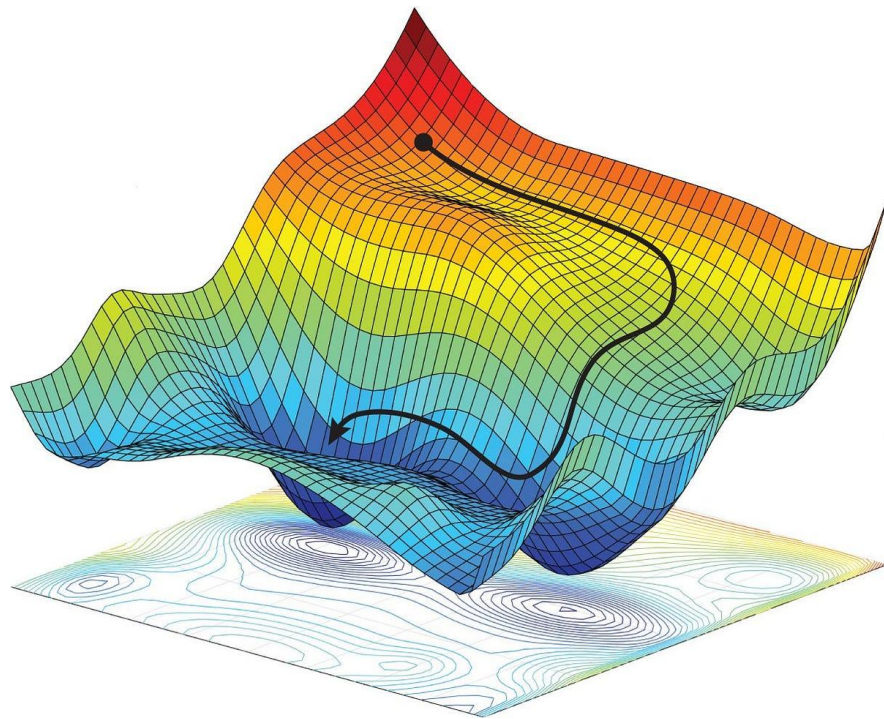
- ❖ Initialize parameters (\mathbf{W} 's and b 's)
- ❖ Repeat until converge:

$$\mathbf{W}_2 = \mathbf{W}_2 - \alpha d\mathbf{W}_2$$

$$b_2 = b_2 - \alpha db_2$$

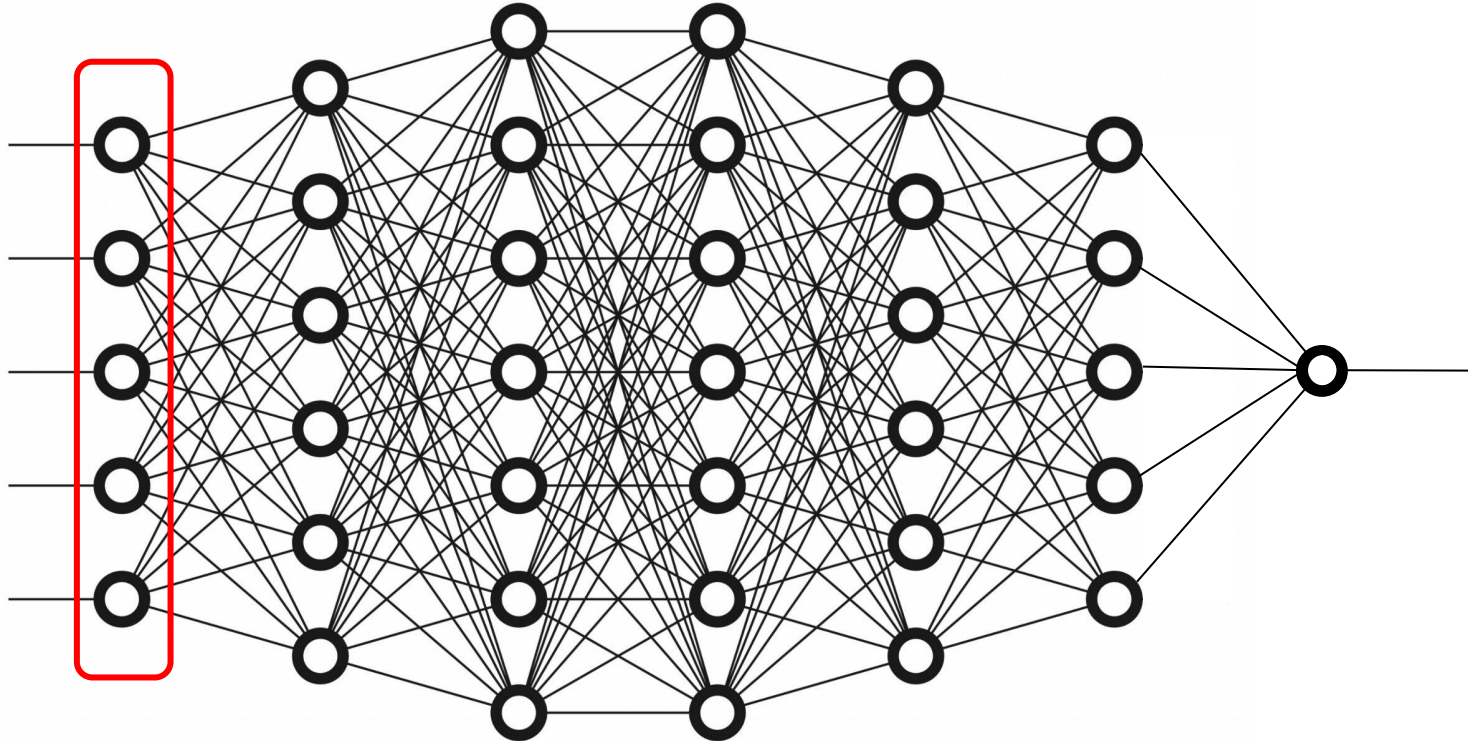
$$\mathbf{W}_1 = \mathbf{W}_1 - \alpha d\mathbf{W}_1$$

$$b_1 = b_1 - \alpha db_1$$



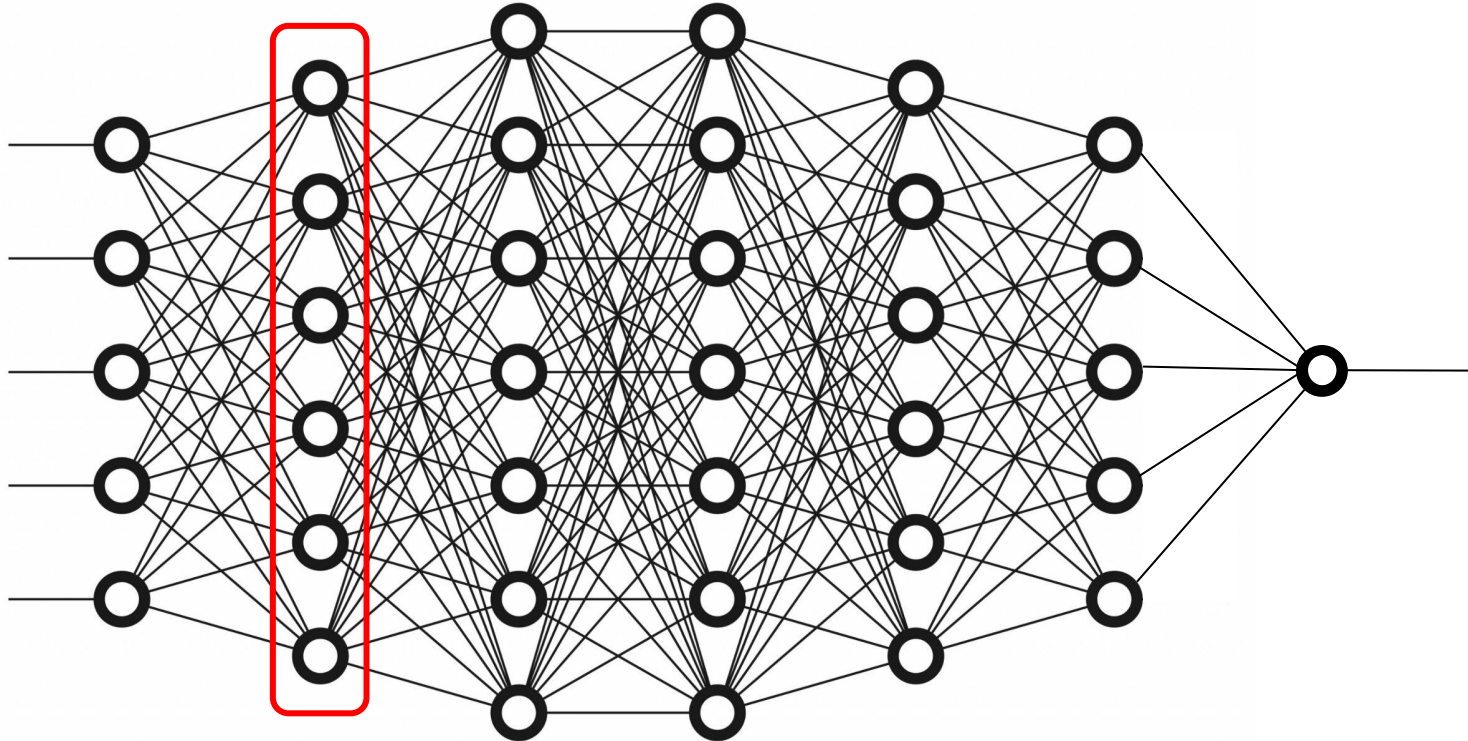
Deep Networks

Forward propagation



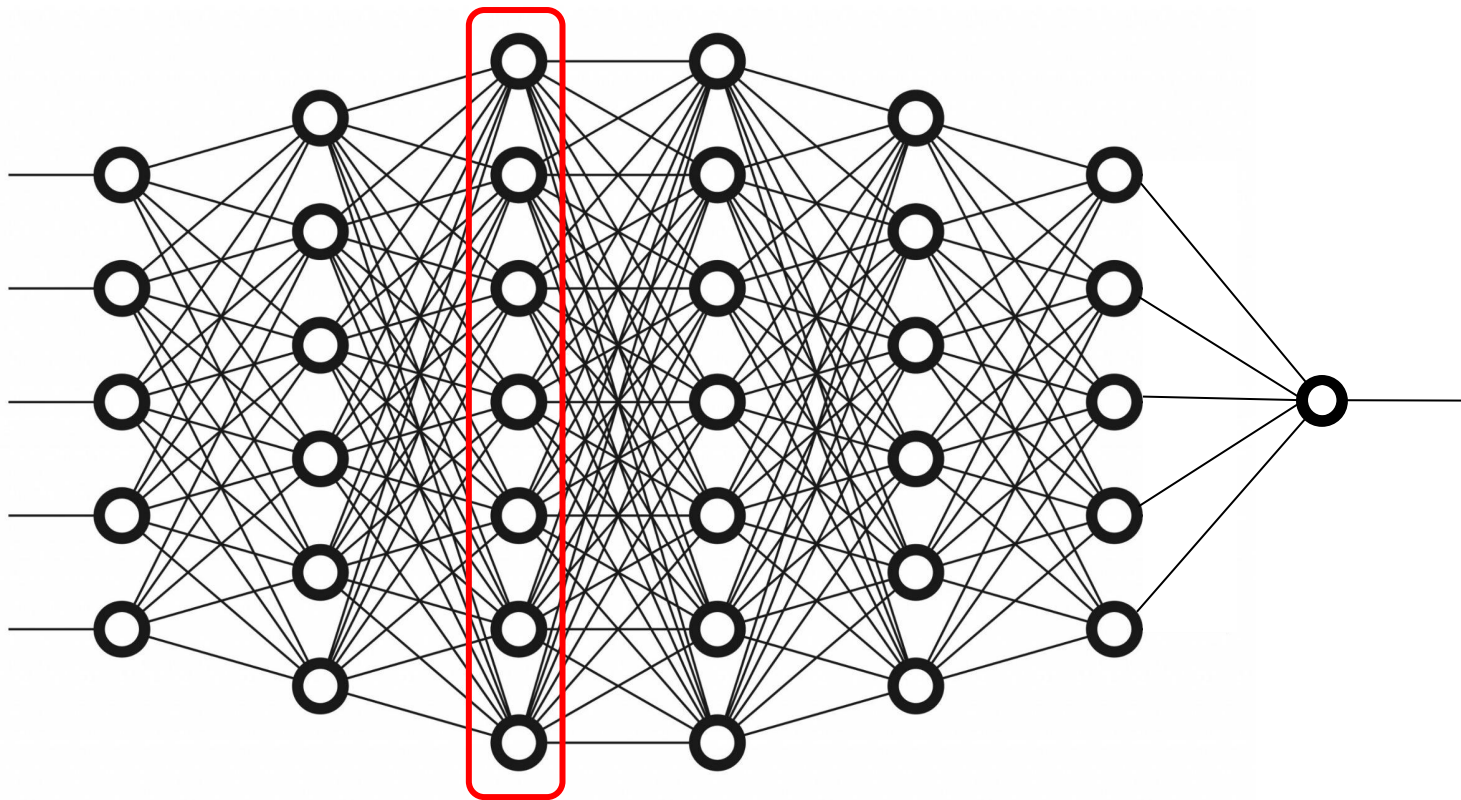
Deep Networks

Forward propagation



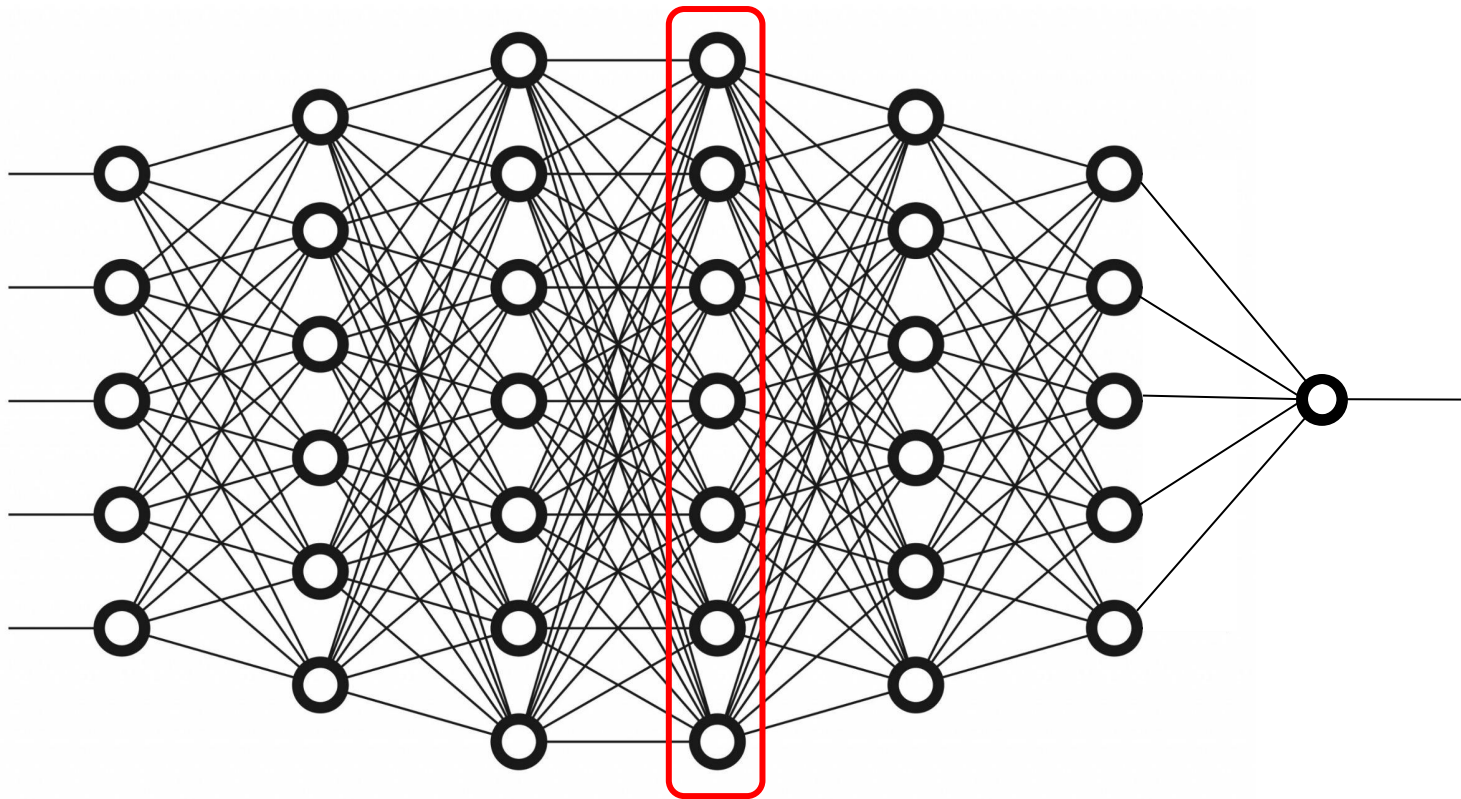
Deep Networks

Forward propagation



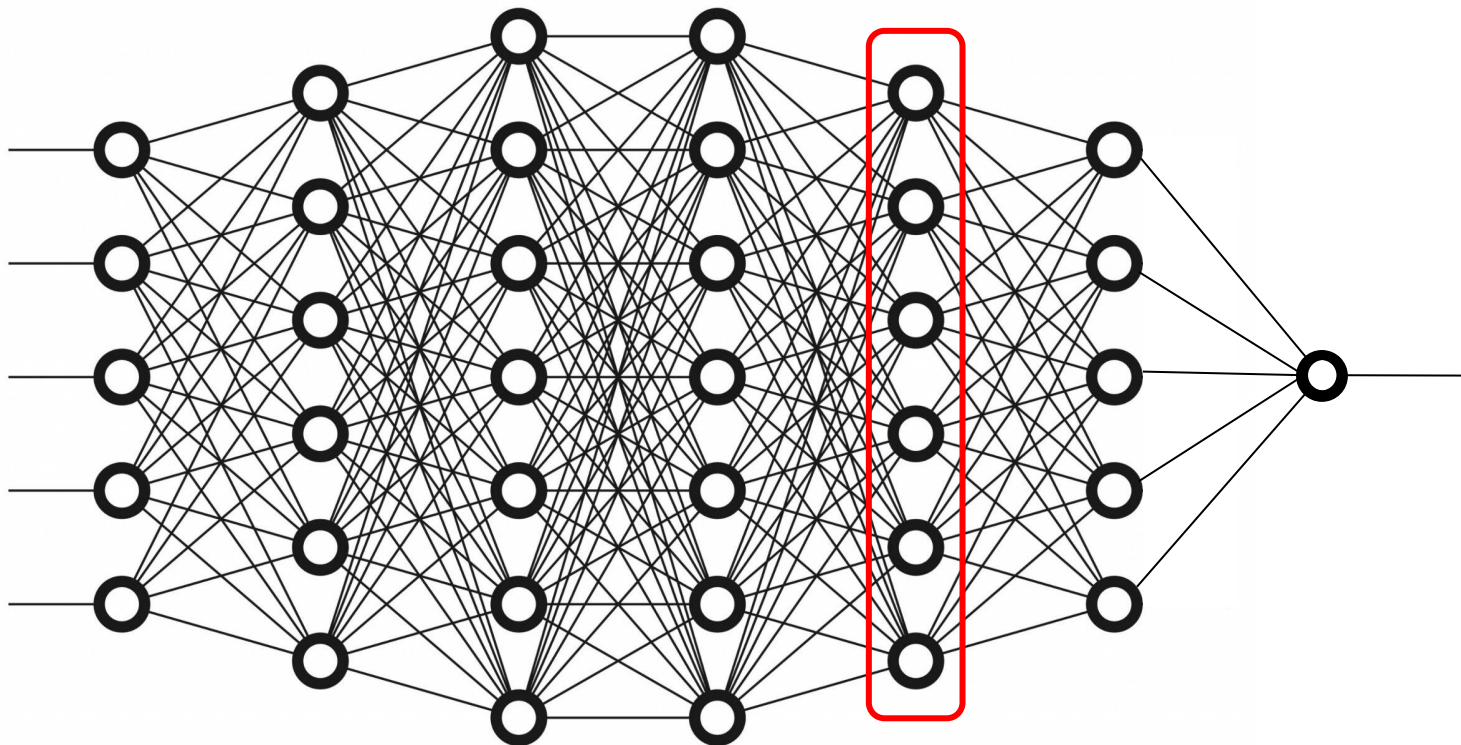
Deep Networks

Forward propagation



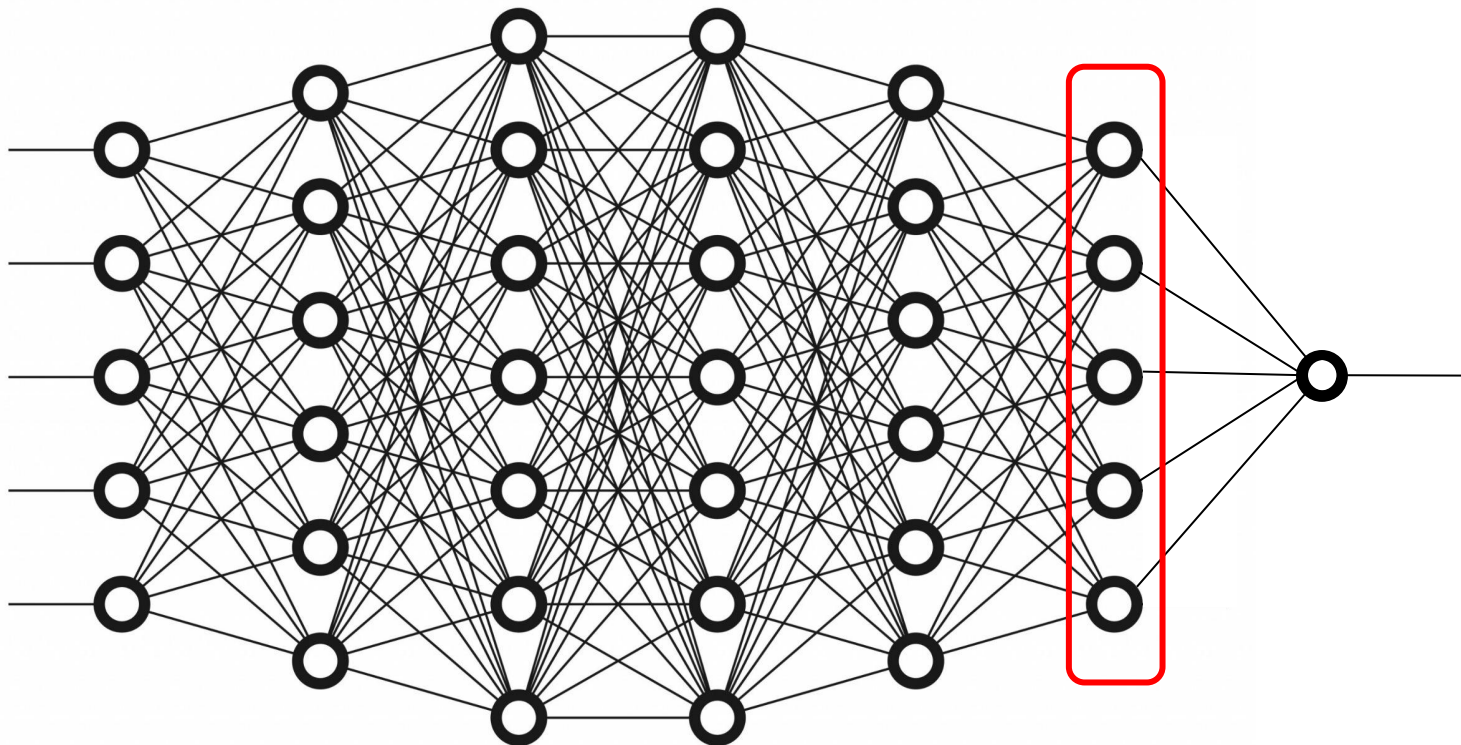
Deep Networks

Forward propagation



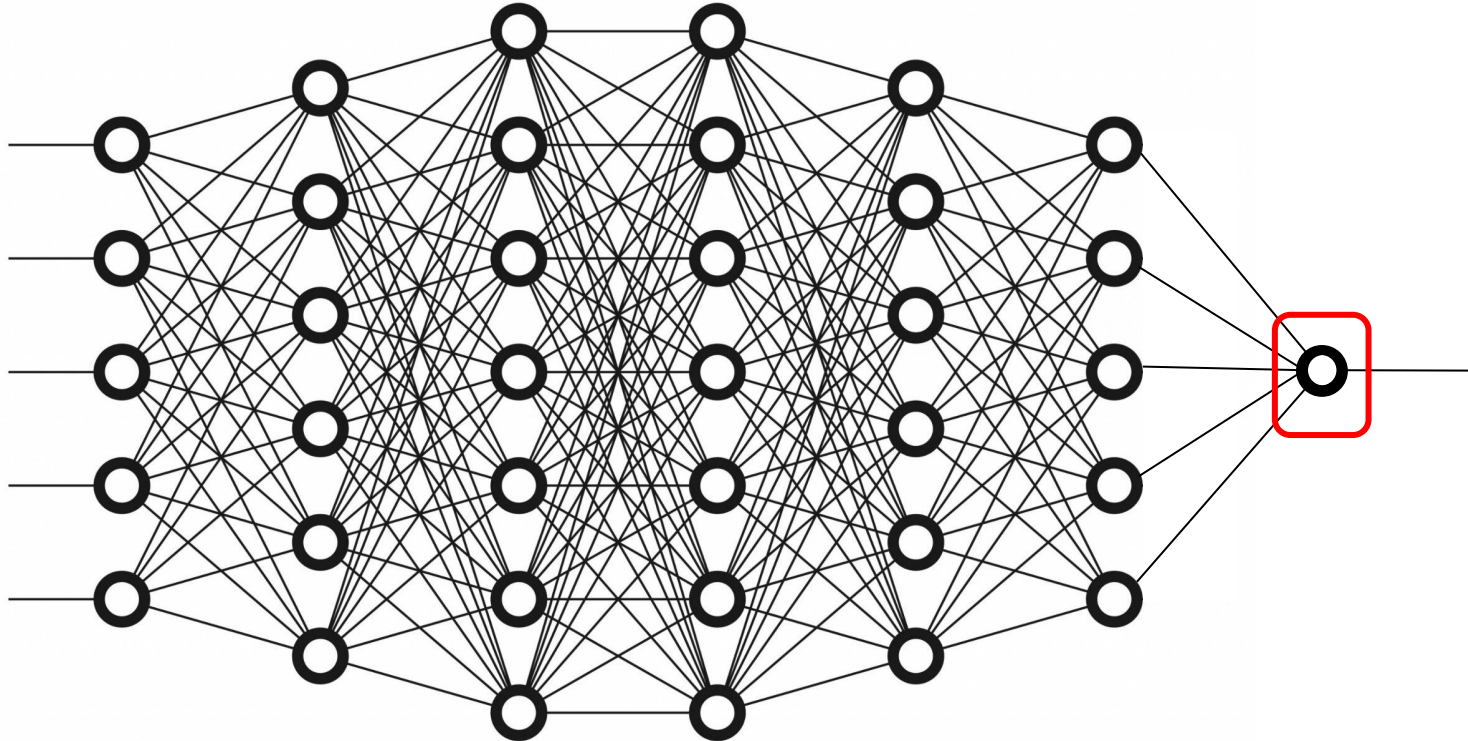
Deep Networks

Forward propagation



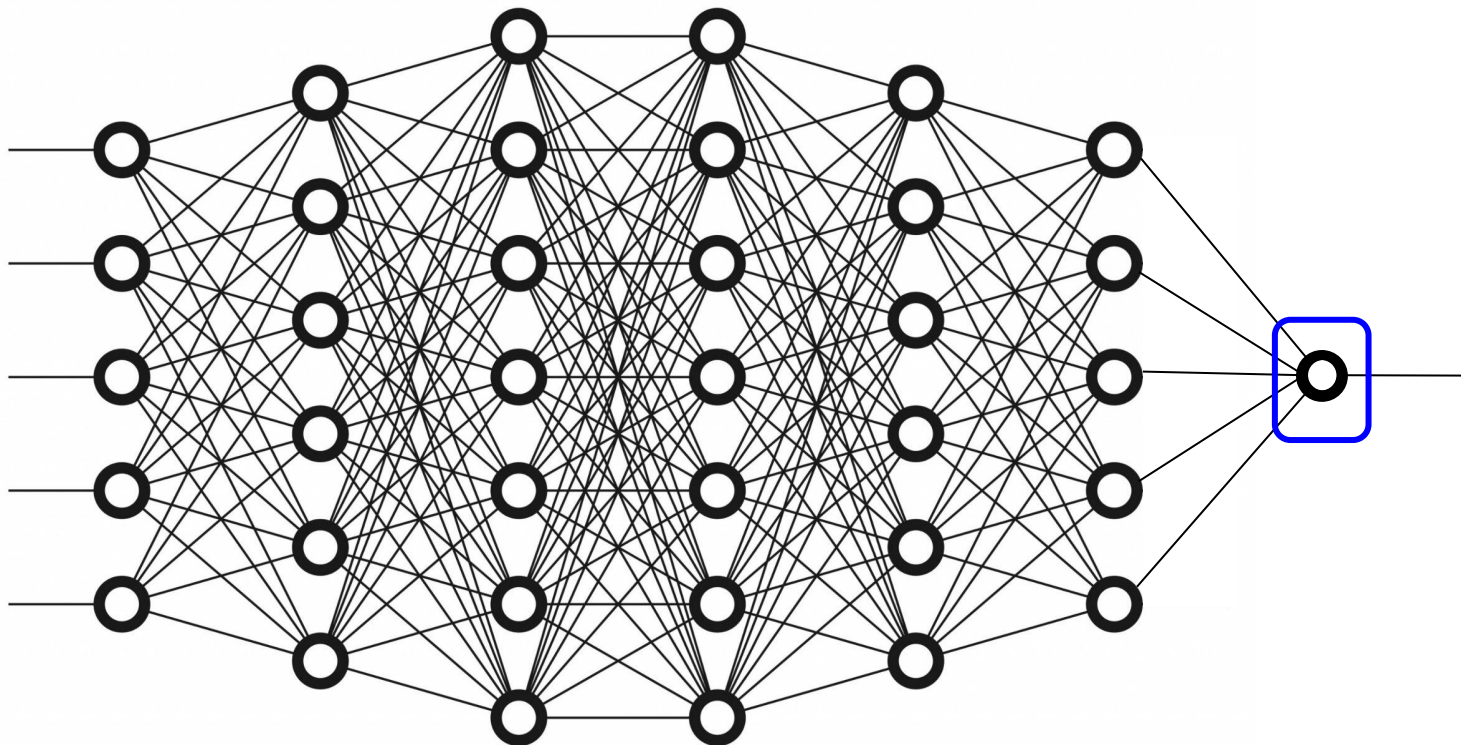
Deep Networks

Forward propagation



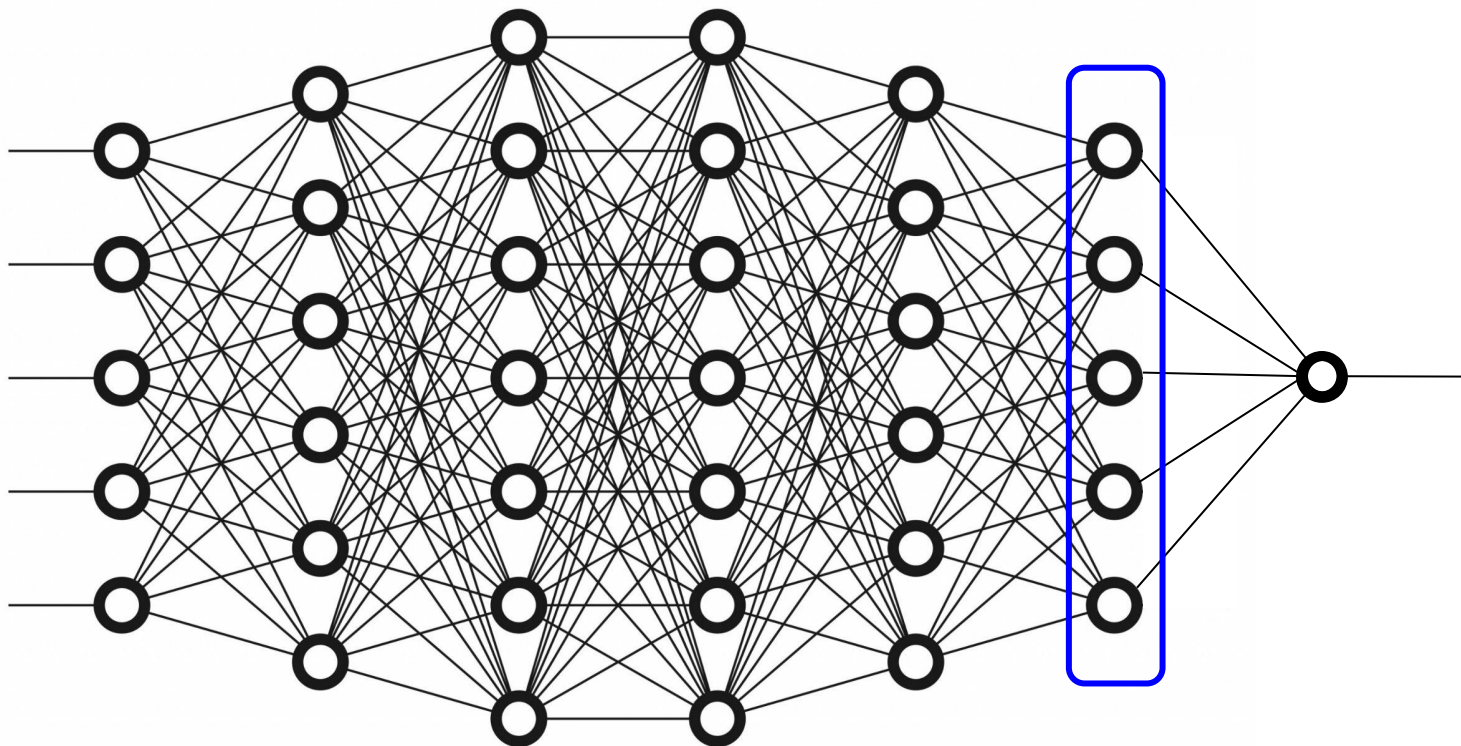
Deep Networks

Backward propagation



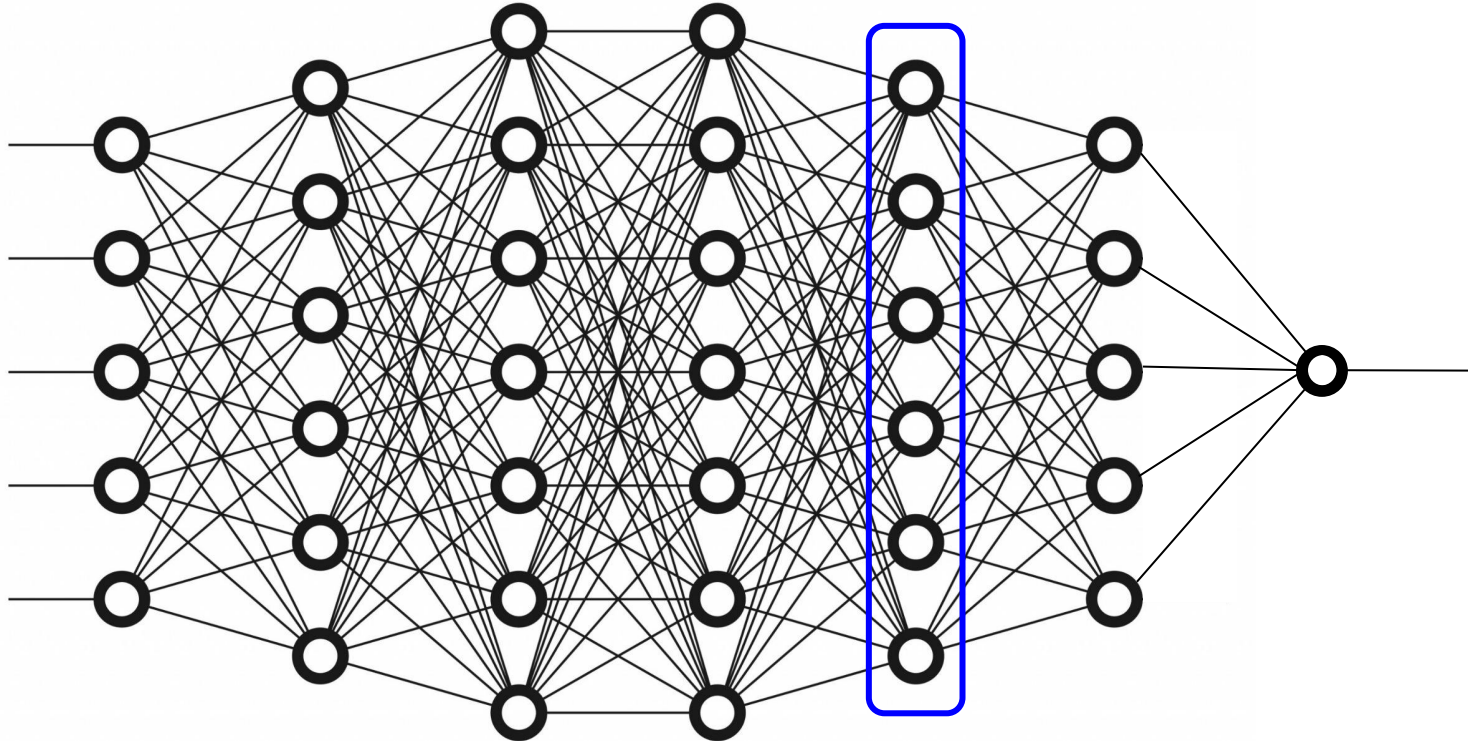
Deep Networks

Backward propagation



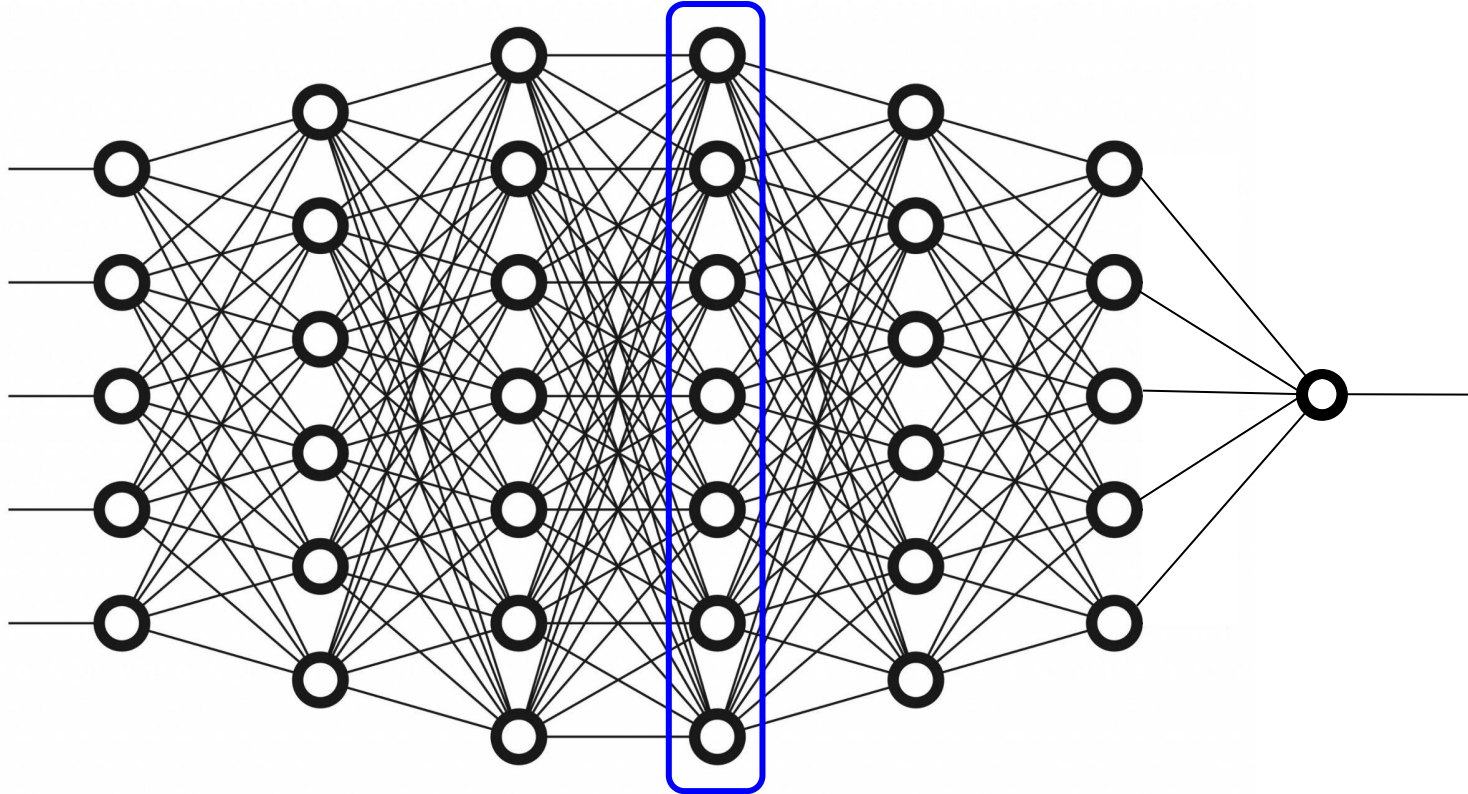
Deep Networks

Backward propagation



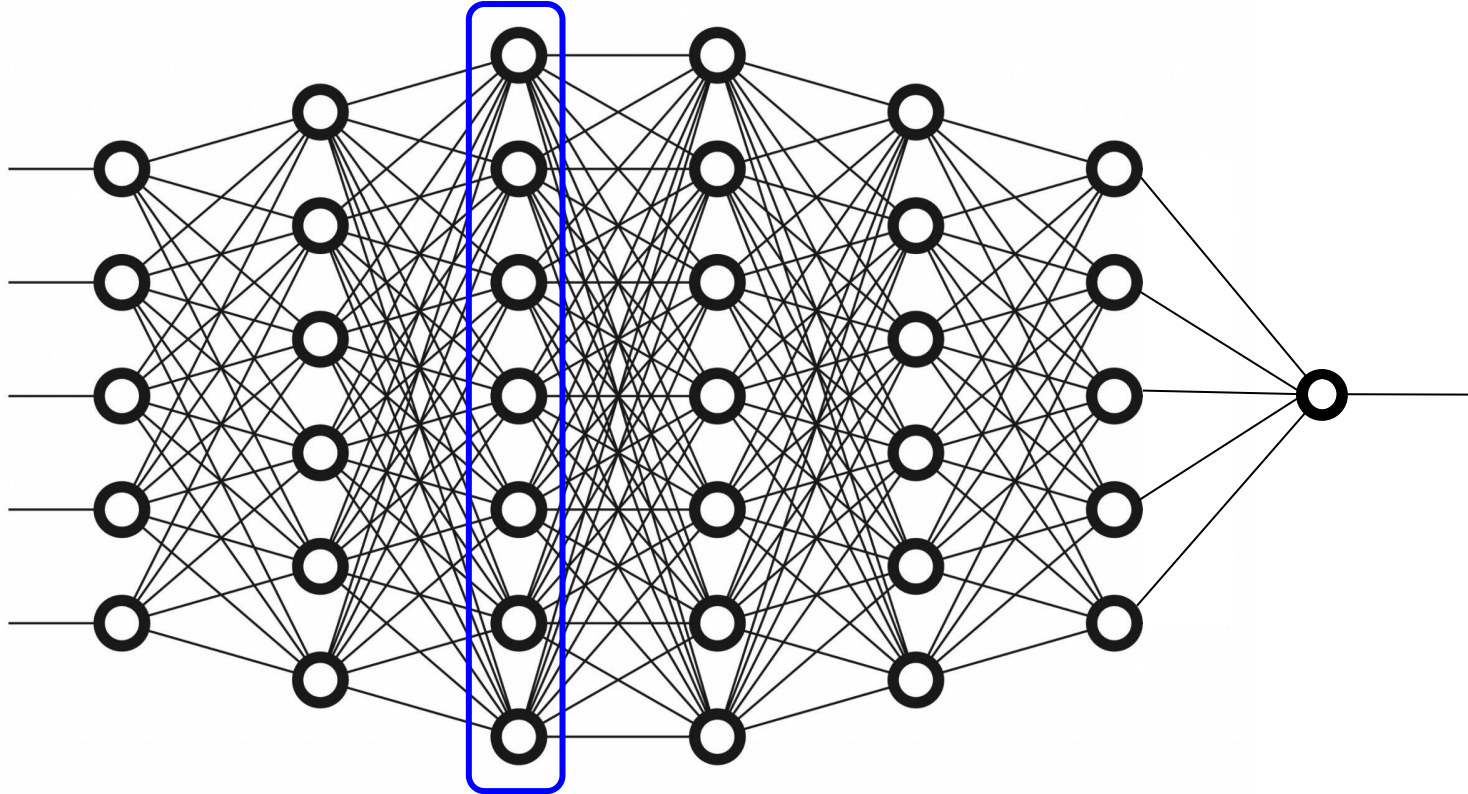
Deep Networks

Backward propagation



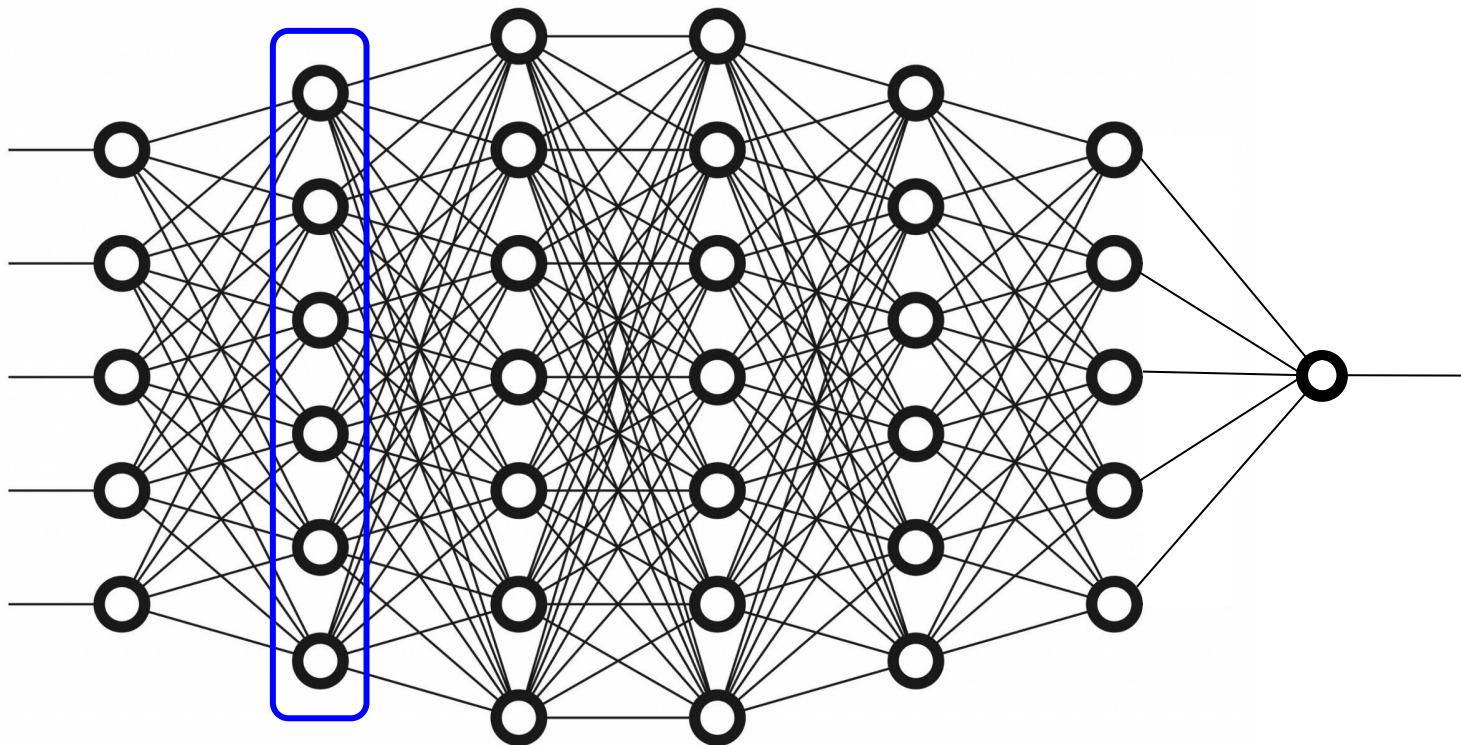
Deep Networks

Backward propagation



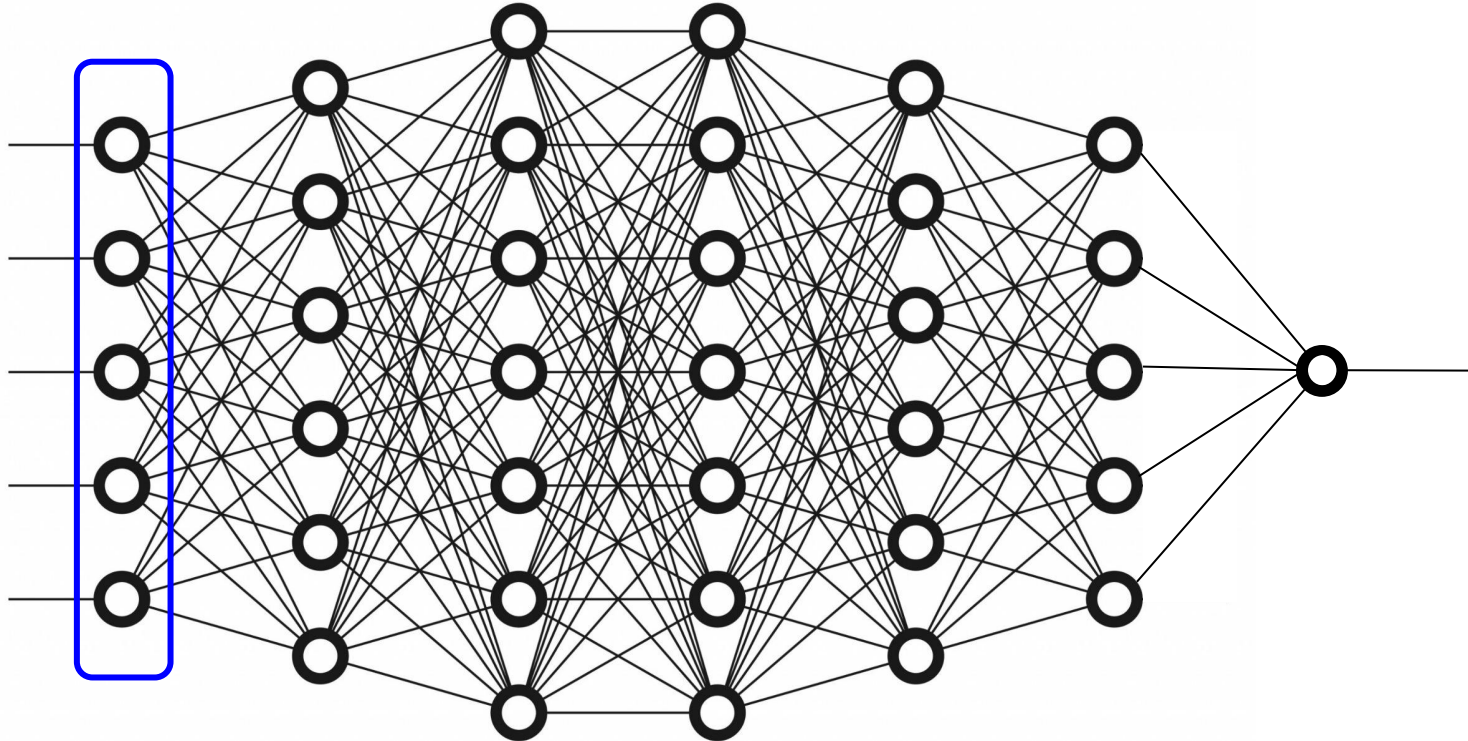
Deep Networks

Backward propagation



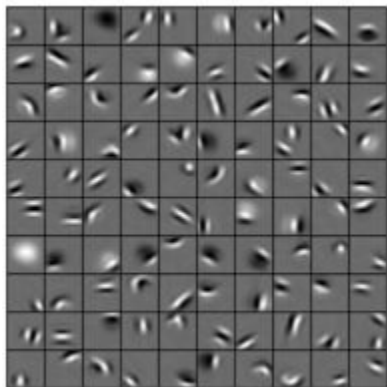
Deep Networks

Backward propagation



What do layers learn

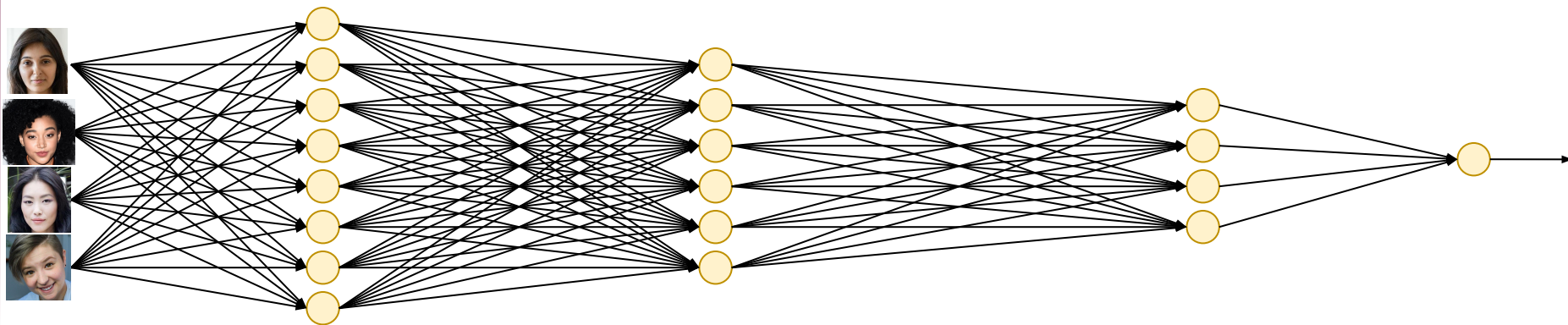
Low-Level



Mid-Level

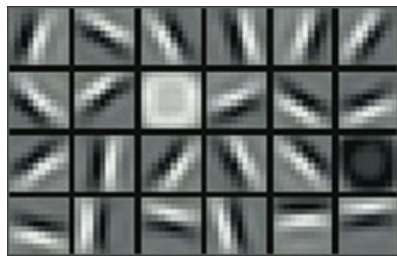


High-Level



What do layers learn

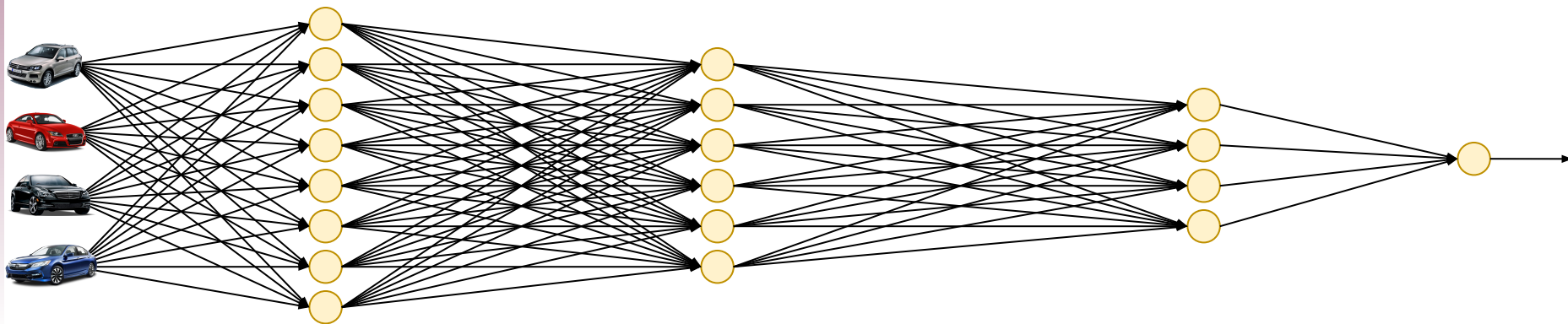
Low-Level



Mid-Level



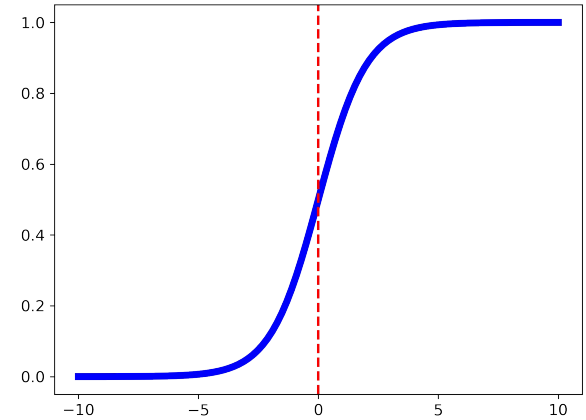
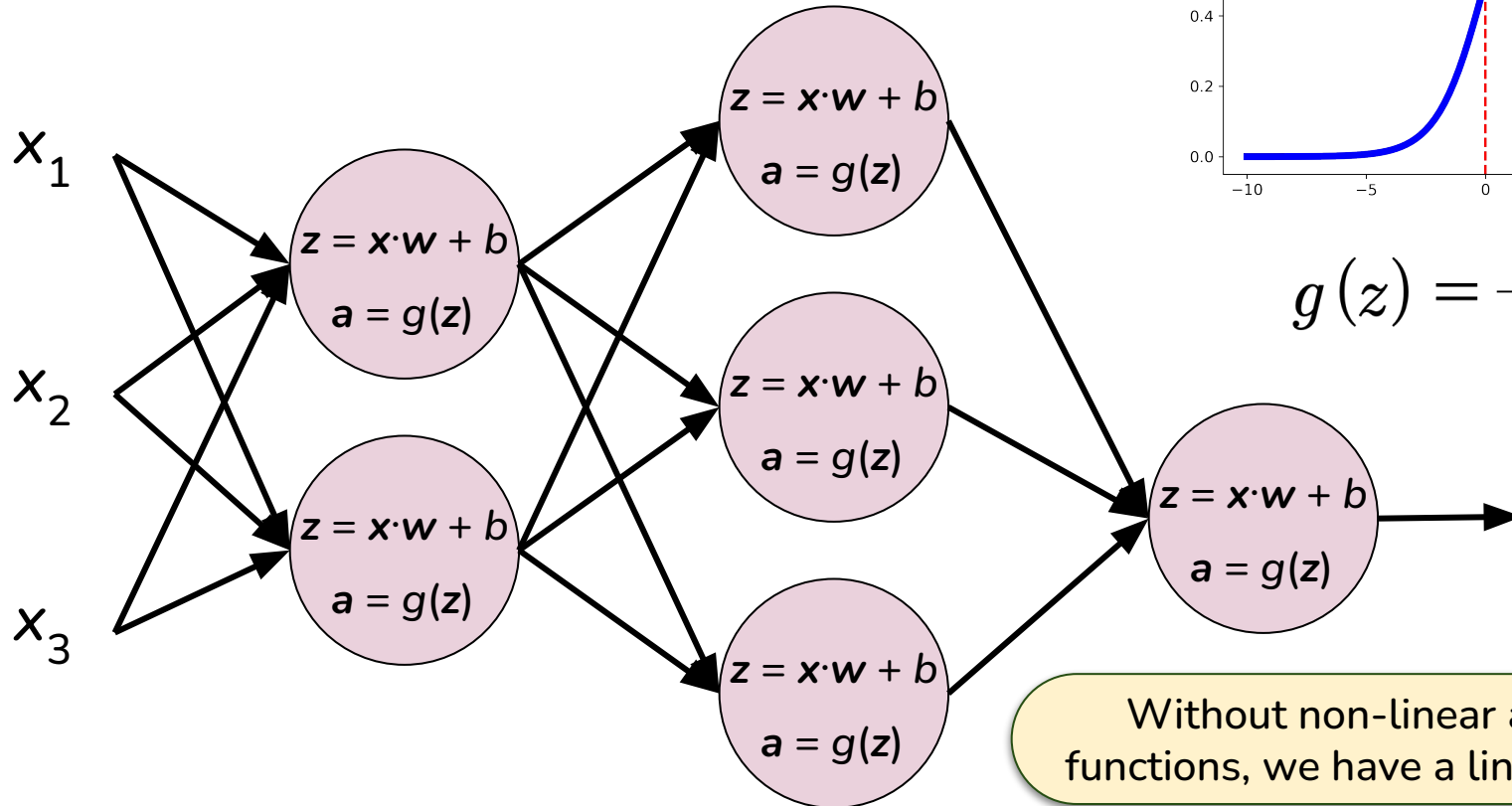
High-Level



Extensions

- ❖ Activation functions
- ❖ Classification and regression
- ❖ Hyperparameter tuning

Activation Functions



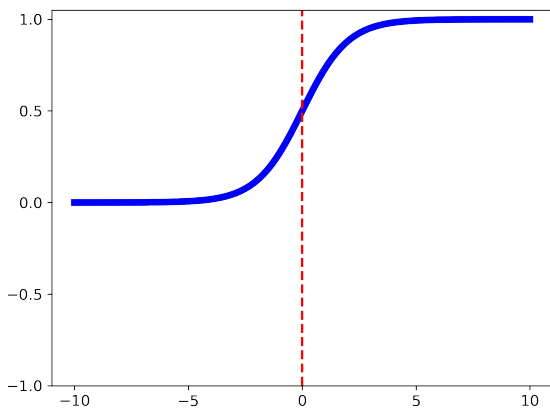
$$g(z) = \frac{1}{1 + e^{-z}}$$

Without non-linear activation functions, we have a linear classifier

Activation Functions

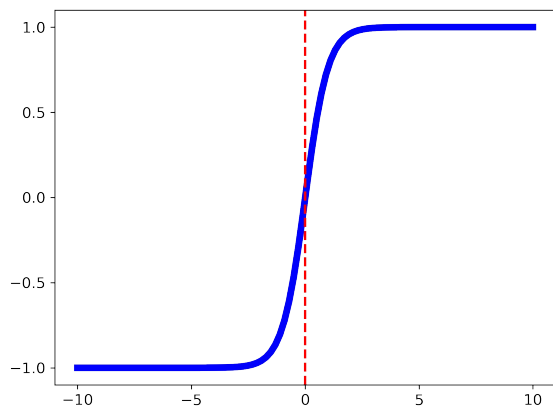
ReLU is most common
for *hidden* layers

sigmoid



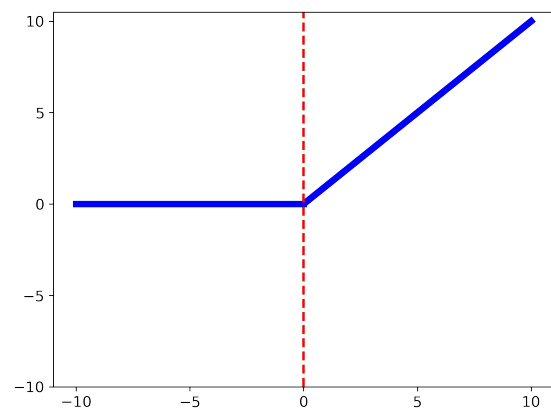
$$\text{sigmoid}(z) = \frac{1}{1 + e^{-z}}$$

tanh
(hyperbolic tangent)



$$\text{tanh}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

ReLU
(Rectified Linear Unit)



$$\text{relu}(z) = \max(0, z)$$

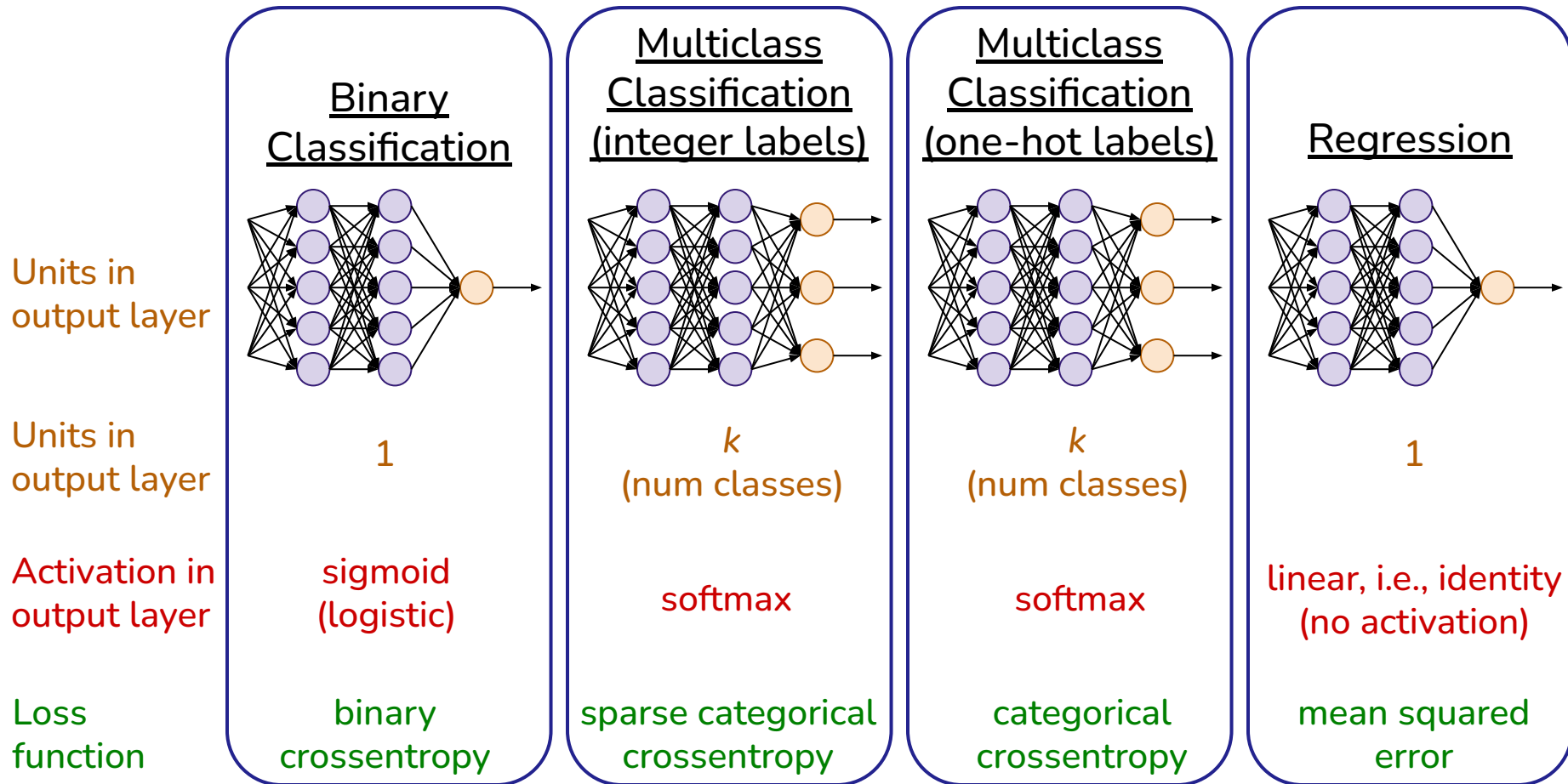
Extensions

❖ Activation functions

❖ Classification and regression

❖ Hyperparameter tuning

Classification and Regression



Extensions

- ❖ Activation functions
- ❖ Classification and regression
- ❖ Hyperparameter tuning

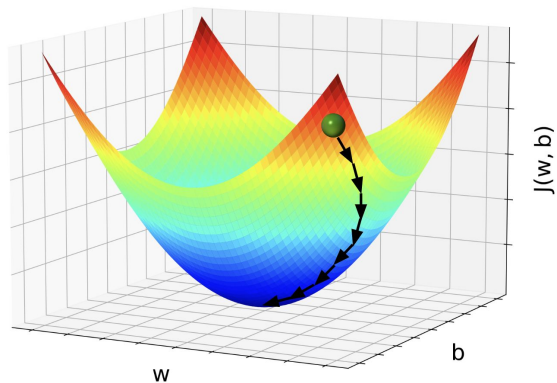
Hyperparameter Tuning

Using *validation* data!

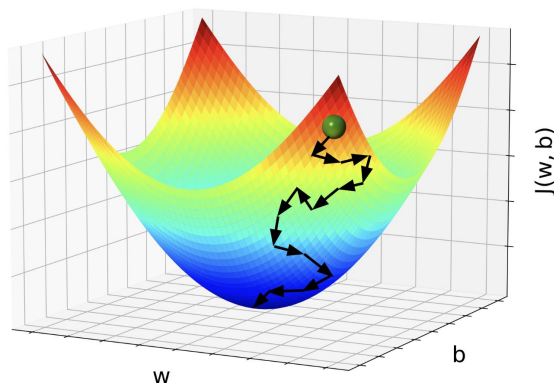
- ❖ Architecture, i.e., number of layers and units per layer
- ❖ Activation functions (ReLU for hidden layers)
- ❖ Batch size (32, 64, 128, 256, 512)

Gradient Descent

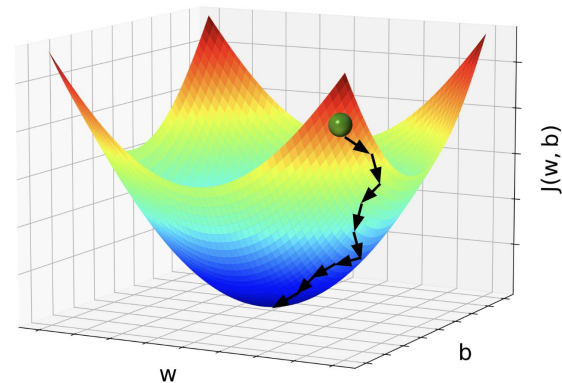
Batch
(batch size = m)



Stochastic
(batch size = 1)



Mini-Batch
(batch size = 64)

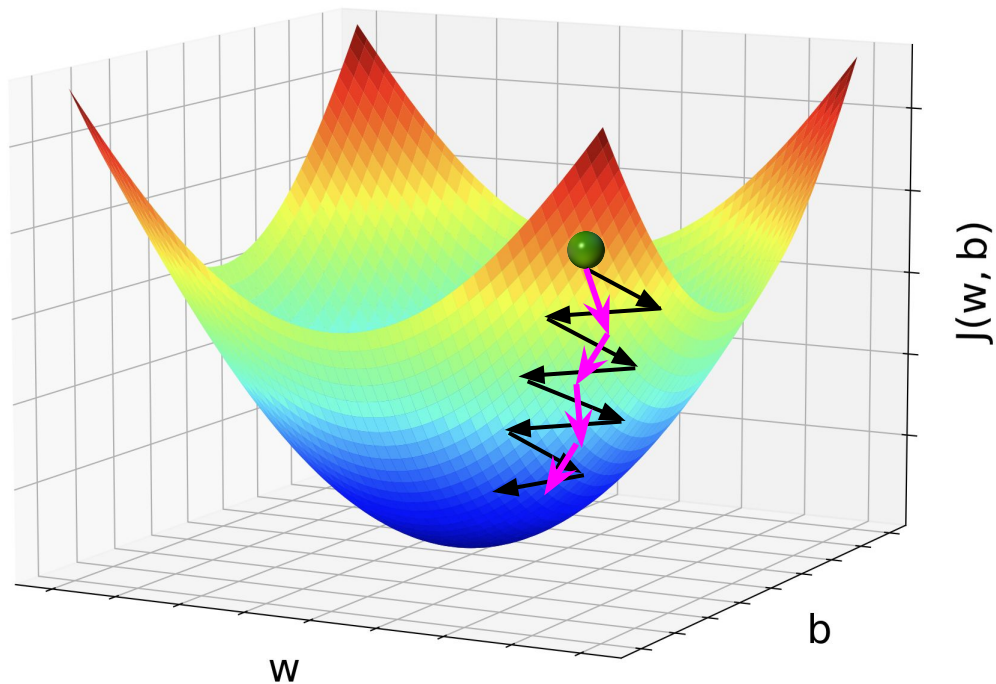
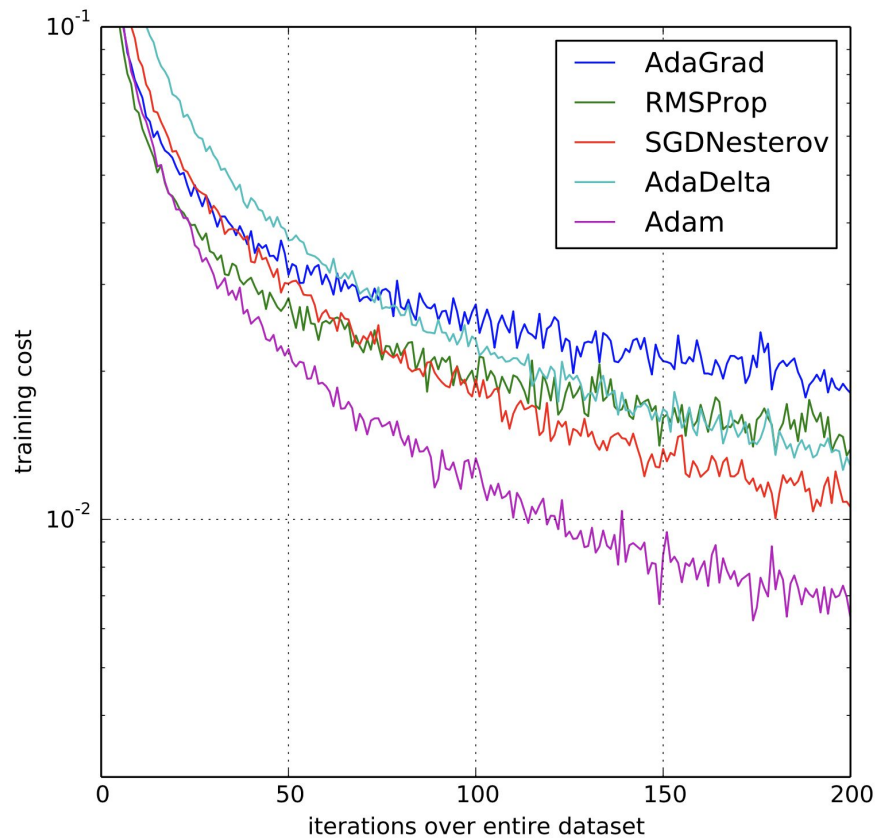


Hyperparameter Tuning

Using *validation* data!

- ❖ Architecture, i.e., number of layers and units per layer
- ❖ Activation functions (ReLU for hidden layers)
- ❖ Batch size (32, 64, 128, 256, 512)
- ❖ Learning rate, α . Decrease over time.
- ❖ Iterations of gradient descent (convergence and **max_iter**)
- ❖ Gradient descent (Adam: adaptive moment estimation)

Improvements on Gradient Descent



Hyperparameter Tuning

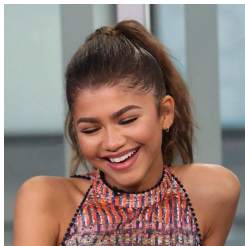
Using *validation* data!

- ❖ Architecture, i.e., number of layers and units per layer
- ❖ Activation functions (ReLU for hidden layers)
- ❖ Batch size (32, 64, 128, 256, 512)
- ❖ Learning rate, α . Decrease over time.
- ❖ Iterations of gradient descent (convergence and `max_iter`)
- ❖ Gradient descent (Adam: adaptive moment estimation)
- ❖ Regularization (L2, dropout, early stopping)

Overfitting

- ❖ **Overfitting** is one of the most common problems in ML
- ❖ Model learns properties specific to the training data that don't generalize to new (testing) data
- ❖ Performance is much better on training data than on testing data

Z



TS

Regularization - L_2

Smaller values for the parameters w lead to more generalizable models and are less prone to overfitting

To incentivize small values for w , modify the cost function so that it:

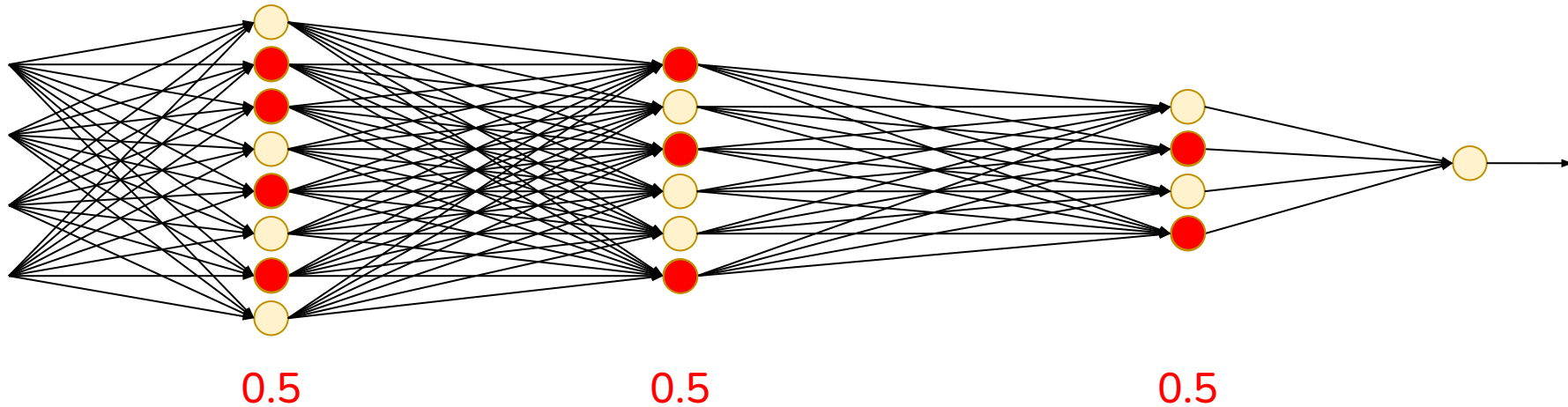
- 1) Fits the training data well
- 2) Penalizes large values for w

and

λ is
regularization
parameter

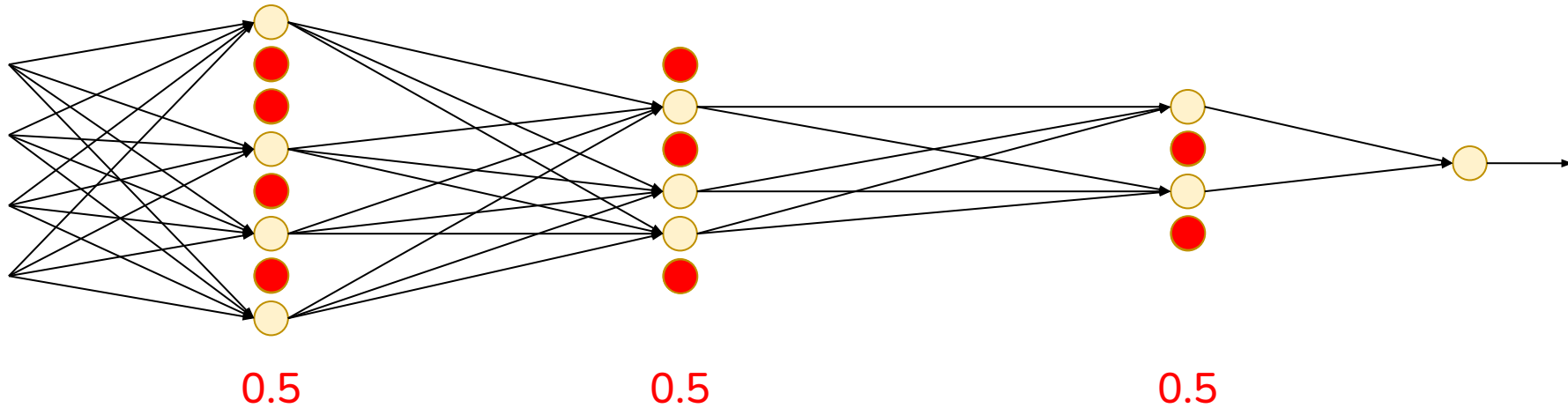
$$J = \frac{1}{m} \sum_{i=1}^m L + \frac{\lambda}{2m} \sum \|w\|^2$$

Regularization - Dropout



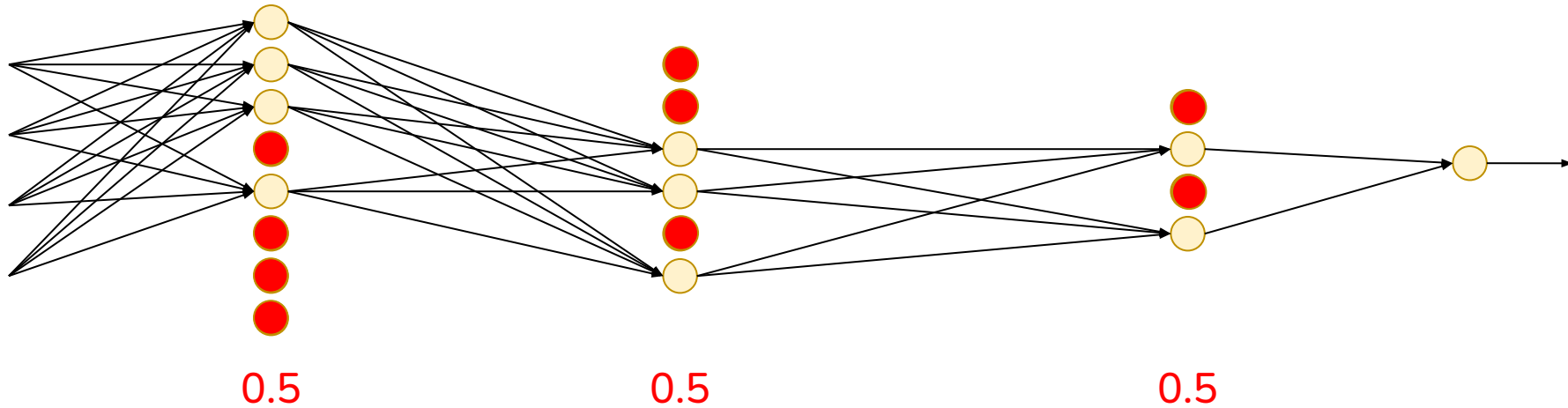
Randomly choose units to remove from network
each time parameters are updated

Regularization - Dropout



Randomly choose units to remove from network
each time parameters are updated

Regularization - Dropout



Randomly choose units to remove from network each time parameters are updated

Regularization - Early Stopping

